TMA4267 Linear Statistical Models V2017 [L7] Part 2: Linear regression [F p73-86] Model definition [F3.1], Parameters and residuals [F3.1.1], Model check [F3.1.2]

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: February 7, 2017

Part 2: Linear regression

Part 2: Linear regression

Fahrmeir et al (2013): Regression. Chapter 3.1, 3.2, 3.4 and required parts of 3.5 and Appendix B.

Part 3: Hypothesis testing and analysis of variance

- Fahrmeir et al (2013): Regression. Chapter 3.3 and required parts of 3.5 and Appendix B.
- Härdle et al (2015): Applied Multivariate Statistical Analysis. Chapter 8.1.1. (ANOVA).
- A short note on multiple testing (to be written).

File TMA4267Part2and3.pdf available from course www-page.

Age-predicted maximal heart rate in healthy subjects: The HUNT Fitness Study

B. M. Nes, I. Janszky, U. Wisløff, A. Støylen, T. Karlsen (2012) in Scandinavian Journal of Medicine and Science in Sports.

- HRmax describes the highest heart rate achieved by a subject exercising to exhaustion and is verified by a plateau of heart rate despite increasing workload. In the literature, HRmax commonly refers to the peak heart rate at termination of a graded maximal exercise test.
- However, in clinical settings, a maximal exercise test is not always feasible and there is a need to predict HRmax from age prior to testing to be able to adequately assess heart rate response and relative intensity of effort at submaximal levels.

Age-predicted maximal heart rate in healthy subjects: The HUNT Fitness Study

- HRmax at a given age is frequently estimated by the "220 age" formula.
- The aim of the present study was to develop a new prediction formula for HRmax through analysis of HRmax measured at VO2peak in a diverse population of 4635 healthy subjects and compare this formula with three commonly used prediction formulas. Furthermore, we wanted to investigate the relationship between HRmax and gender, physical activity status, BMI, and objectively measured aerobic fitness.

Age-predicted maximal heart rate in healthy subjects: The HUNT Fitness Study - Statistical procedures

- Only subjects that fulfilled the criteria of a maximal test, with registered maximal heart rate (HRmax), were included in the analysis (n = 3320).
- General linear modeling was used to determine the effect of age on HRmax. HRmax was entered as the dependent variable and age as the independent variable. Nonlinearity of the relationship between age and HRmax was investigated by including polynomial terms to the regression model.
- In a subsequent analysis, the effects of gender, BMI, physical activity status, and maximal oxygen uptake were examined by entering these factors as independent variables in addition to age. In further subsequent models, interaction terms were included as well to assess effect modification.
- The continuous variables were checked for normality, homogeneity of variances, and heteroscedasticity of the residuals.



Nes et al (2012): Age-predicted maximal heart rate in healthy subjects: The HUNT Fitness Study. n = 3320 individuals.

Munich Rent Index data set

described in Fahrmeir et al (2013) on pages 19-20.

- > library("gamlss.data")
- > ds=rent99
- > dim(ds)
- [1] 3082 9
- > colnames(ds)

[1] "rent" "rentsqm" "area" "yearc" "location" "bath"
[7] "kitchen" "cheating" "district"

> summary(ds)

rent	rentsqm	area	yearc
Min. : 40.51	Min. : 0.4158	Min. : 20.00	Min. :1918
1st Qu.: 322.03	1st Qu.: 5.2610	1st Qu.: 51.00	1st Qu.:1939
Median : 426.97	Median : 6.9802	Median : 65.00	Median :1959
Mean : 459.44	Mean : 7.1113	Mean : 67.37	Mean :1956
3rd Qu.: 559.36	3rd Qu.: 8.8408	3rd Qu.: 81.00	3rd Qu.:1972
Max. :1843.38	Max. :17.7216	Max. :160.00	Max. :1997
location bath	kitchen cheating	district	
1:1794 0:2891	0:2951 0: 321	Min. : 113	
2:1210 1: 191	1: 131 1:2761	1st Qu.: 561	
3: 78		Median :1025	
		Mean :1170	
		3rd Qu.:1714	
		Max. :2529	

The classical linear model

The model

$$\mathbf{Y} = \mathbf{X}oldsymbol{eta} + arepsilon$$

is called a classical linear model if the following is true:

1.
$$E(\varepsilon) = 0$$
.

2.
$$\operatorname{Cov}(\varepsilon) = \operatorname{E}(\varepsilon \varepsilon^{T}) = \sigma^{2} I$$
.

3. The design matrix has full rank, rank(X) = k + 1 = p. The classical *normal* linear regression model is obtained if

additionally

4. $\varepsilon \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$

holds. For random covariates these assumptions are to be understood conditionally on X.

Conditional mean and covariance

If we believe that the vector with elements Y and X are multivariate normal $N_{k+1}(\mu, \Sigma)$ we may look at the partition

$$\begin{pmatrix} Y \\ \boldsymbol{X} \end{pmatrix} \sim N_{k+1} \left(\begin{pmatrix} \mu_{Y} \\ \boldsymbol{\mu}_{\boldsymbol{X}} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{YY} & \boldsymbol{\Sigma}_{YX} \\ \boldsymbol{\Sigma}_{XY} & \boldsymbol{\Sigma}_{XX} \end{pmatrix} \right)$$

The conditional distributions of the components are (multivariate) normal, with conditional mean and variance of $Y \mid \mathbf{X} = \mathbf{x}$ are

$$E(Y \mid \boldsymbol{X} = \boldsymbol{x}) = \mu_Y + \boldsymbol{\Sigma}_{YX} \boldsymbol{\Sigma}_{XX}^{-1} (\boldsymbol{x} - \boldsymbol{\mu}_X)$$
$$Var(Y \mid \boldsymbol{X} = \boldsymbol{x}) = \boldsymbol{\Sigma}_Y - \boldsymbol{\Sigma}_{YX} \boldsymbol{\Sigma}_{XX}^{-1} \boldsymbol{\Sigma}_{XY}$$

Observe: mean is linear in x and variance independent of x.

Model assumptions for the classical linear model [F:3.1.2]

What are our model assumptions, how can we spot violations and what can we do to amend the violations.

- 1. Linearity of covariates: $\mathbf{Y} = \mathbf{X}\mathbf{\beta} + \mathbf{\varepsilon}$
- 2. Homoscedastic error variance: $Cov(\varepsilon) = \sigma^2 I$.
- **3.** Uncorrelated errors: $Cov(\varepsilon_i, \varepsilon_j) = 0$.
- 4. Additivity of errors: $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

We mainly use plots to assess this (more on model fit in F:3.4 Model choice and variable seletion)

- Covariate vs response (for each covariate)
- Covariate vs error (when we have simulated data and know the truth)
- Covariate vs residual (estimated error),
- Predicted response vs residual (to be popular later).

Linearity of covariates: Covariate vs. response

Munich Rent Index: area vs rentsqm



Linearity of covariates: Covariate vs. residual (residual plot) Munich Rent Index: area vs residual



ds\$area

Linearity of covariates: Transformed covariate vs. response Munich Rent Index: 1/area vs rentsqm



1/ds\$area

Linearity of covariates: Transformed covariate vs. residual (residual plot)

Munich Rent Index: 1/area vs residual



13 / 20

3.2 Modeling Nonlinear Covariate Effects Through Variable Transformation

If the continuous covariate z has an approximately nonlinear effect $\beta_1 f(z)$ with known transformation f, then the model

 $y_i = \beta_0 + \beta_1 f(z_i) + \ldots + \varepsilon_i$

can be transformed into the linear regression model

$$y_i = \beta_0 + \beta_1 x_i + \ldots + \varepsilon_i,$$

where $x_i = f(z_i) - \overline{f}$. By subtracting

$$\bar{f} = \frac{1}{n} \sum_{i=1}^{n} f(z_i),$$

the estimated effect $\hat{\beta}_1 x$ is automatically centered around zero. The estimated curve is best interpreted by plotting $\hat{\beta}_1 x$ against z (instead of x).

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.94)

3.3 Modeling Nonlinear Covariate Effects Through Polynomials

If the continuous covariate z has an approximately polynomial effect $\beta_1 z + \beta_2 z^2 + \ldots + \beta_l z^l$ of degree l, then the model

$$y_i = \beta_0 + \beta_1 z_i + \beta_2 z_i^2 + \ldots + \beta_l z_i^l + \ldots + \varepsilon_i$$

can be transformed into the linear regression model

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_1 x_{i2} + \ldots + \beta_l x_{il} + \ldots + \varepsilon_i,$$

where $x_{i1} = z_i, x_{i2} = z_i^2, \dots, x_{il} = z_i^l$.

The centering (and possibly orthogonalization) of the vectors $\mathbf{x}^{j} = (x_{1j}, \ldots, x_{nj})'$, $j = 1, \ldots, l$, to $\mathbf{x}^{1} - \bar{\mathbf{x}}_{1}, \ldots, \mathbf{x}^{l} - \bar{\mathbf{x}}_{l}$ with the mean vector $\bar{\mathbf{x}}_{j} = (\bar{x}_{j}, \ldots, \bar{x}_{j})'$ facilitates interpretation of the estimated effects. A graphical illustration of the estimated polynomial is a useful way to interpret the estimated effect of z.

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.95)

Homoscedastic errors

```
n=1000
x=seq(-3,3,length=n)
beta 0 = -1
beta1=2
xbeta=beta0+beta1*x
sigma=1
e1=rnorm(n,mean=0,sd=sigma)
y1=xbeta+e1
ehat1=residuals(lm(y1~x))
plot(x,y1,pch=20)
abline(beta0,beta1,col=1)
plot(x,e1,pch=20)
abline(h=0,col=2)
```

Heteroscedastic errors

```
sigma=(0.1+0.3*(x+3))^2
e2=rnorm(n,0,sd=sigma)
y2=xbeta+e2
ehat2=residuals(lm(y2~x))
plot(x,y2,pch=20)
abline(beta0,beta1,col=2)
plot(x,e2,pch=20)
abline(h=0,col=2)
```

Homo- and heteroscedastic errors



Top: homoscedastic errors. Bottom: heteroscedastic errors. Right: x vs y. Left: x vs error. Example from Fahrmeir et al (2013): Regression. Springer. (p.79). R code from TMA4267 lectures tab.

Homoscedastic errors?



Today

- Normal linear model: implication for Y.
- Model parameters β , σ^2 , parameter estimators $\hat{\beta}$, $\hat{\sigma}^2$, residuals $\hat{\varepsilon} = Y \mathbf{X}\hat{\beta}$.
- Model assumptions.
- Next: covariates- how to include in linear regression, and then parameter estimation.

PART Q:
LINEAR REGRESSION
Model definition EF31.0]
Y = Variable of primary interest Gresponde dependent
X₁, X₂,.., X_k = regressors, explenetory vanables
independent variables, conjoination
Assumptions:
Y = f(X₁, X₂,..., X_k) + E
supplementic
1) Systematic component is a linear combination
of the covariations
f(X₁, X₂,..., X_k) = Bot Bixi + Fixet+ BuXk
Simple
X =
$$\begin{bmatrix} X_n \\ X_2 \\ X_{k} \end{bmatrix} = \begin{bmatrix} p_1 \\ p_k \end{bmatrix} = f(x) = X^T B$$

p= ket
2) Additiste errors Y = XTB + E
Reptrictive? Haybe 3 brensformations?

Deta end debign metrix
We collect independent data
$$(Y_{i}, Y_{i})$$
 for $i = 1,.., n$
response
 $Y = \begin{bmatrix} Y_{i} \\ Y_{i} \\ Y_{i} \end{bmatrix}$, $E = \begin{bmatrix} E_{i} \\ E_{i} \\ E_{n} \end{bmatrix}$
 $X = \begin{bmatrix} 1 & X_{i1} & X_{i2} & \cdots & X_{ik} \\ \vdots \\ X_{i1} \end{bmatrix} \begin{bmatrix} 1 & E_{i1} \\ E_{i1} \\ E_{i1} \end{bmatrix}$
 $X = \begin{bmatrix} 1 & X_{i1} & X_{i2} & \cdots & X_{ik} \\ \vdots \\ 1 & X_{i1} & X_{i2} & X_{i1} \end{bmatrix}$
 $X = \begin{bmatrix} 1 & X_{i1} & X_{i2} & \cdots & X_{ik} \\ \vdots \\ 1 & X_{i1} & X_{i2} & X_{i1} \end{bmatrix}$
 $A = number of observations$
 $P = number of observations
 $P = number of covernations$
 $P = number of covernations + 1 (intercept)$
 $Q: What can make reach $(Z) = p$
 $E_{X:}: Hunich reat index; n = 3082$
 $Y = reat or reat program $\Rightarrow Y = \begin{bmatrix} E_{26} \\ 0.41 \\ \vdots \end{bmatrix}$$$$

$$B = \begin{bmatrix} 1 & e_{0} & (a_{1}b_{0} & 1 & 0 & 0 & 0 & a_{1b} \\ 1 & e_{1} & e_{1} & 1 & 1 & 1 & 1 \\ 1 & e_{1} & a_{1} & a_{1} & 1 & 1 & 1 \\ 1 & e_{1} & a_{2} & 1 & 1 & 1 & 1 \\ 1 & e_{1} & a_{2} & 1 & 1 & 1 & 1 \\ 1 & e_{1} & a_{2} & 1 & 1 & 1 & 1 \\ 1 & e_{1} & a_{2} & 1 & 1 & 1 & 1 \\ 1 & e_{1} & a_{2} & 1 & 1 & 1 & 1 \\ 1 & e_{1} & a_{2} & 1 & 1 & 1 & 1 \\ 1 & e_{1} & e_{1} & 1 & 1 & 1 & 1 \\ 1 & e_{1} & e_{1} & e_{1} & e_{1} & e_{1} & e_{$$

$$Y \sim Nn$$
 since $Y = X_{3} + E$, and
 $Constant = Nn$
 $E(Y) = E(X_{3} + E) = X_{3} + E(E) = X_{3}$
 $Cov(Y) = Cov(X_{3} + E) = 0 + Cov(E) = 0^{2}I$

 $\gamma \sim N_n(X_\beta, o^2I)$

The coverates X may be regered as rendom vanables, and then the assumptions (1)+(2) are made conditional on X = X, so E(E | X = X) = 0 and $Cov(E | X = X) = T^T$

If we isleed assume that

$$\begin{bmatrix} Y \\ X_n \\ X_n \\ X_n \\ X_n \\ X_n \end{bmatrix} \sim N_{ne_1} \Rightarrow E(Y | X=x) = linear in x$$

 $Vor(Y | X=x) = not dependent on x$

Model paremeters, estimates endresiduals [F 3.1.1] Y= \$\$\$+8, E(E)=0, Car(E)=6 I The model parameter ar B, 02 the unknown px1 (LS) We will develop potructors: B= (XTX)-1XTY by least squares and meximum likelihood. G² = h-p (Y-Xβ)^T(Y-Xβ) by restricted meximum like likes (REAL) Further: Y is a rendom vector with mean $X\beta$, end estimator for $E(Y) = X\beta$ is $\hat{Y} = X\beta$. the error E is a rendom rector with E(E) = 0 = not Or (E) = 5° I, but E is not observed. L'unham Y = X & + E ~ unobserved Observed Observed Our best guess for the error is the residual vector (É, e) E= Y-Y=Y-XB So, the residuals can be calculated, and we may think of the residuals as predictions of the errors 5

we may use a so-called general linear model, with weighted [east squares (lec Ex 3. P4), but Z is in general ununown.

 \mathcal{A}

TMA4267 Linear Statistical Models V2017 (L8) Part 2: Linear regression: Modelling the effects of covariates [F:3.1.3] Parameter estimation: Estimator for β [F:3.2.1]

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: February 10, 2017

The classical linear model

The model

$$\mathbf{Y} = \mathbf{X}oldsymbol{eta} + arepsilon$$

is called a classical linear model if the following is true:

1.
$$E(\varepsilon) = 0$$
.

2.
$$\operatorname{Cov}(\varepsilon) = \operatorname{E}(\varepsilon \varepsilon^{T}) = \sigma^{2} I$$
.

3. The design matrix has full rank $rank(\mathbf{X}) = k + 1 = p$. The classical *normal* linear regression model is obtained if additionally

4. $\varepsilon \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$

holds. For random covariates these assumptions are to be understood conditionally on X.

Model assumptions for the classical linear model [F:3.1.2]

What are our model assumptions, how can we spot violations and what can we do to amend the violations.

1. Linearity of covariates: $\mathbf{Y} = \mathbf{X}\mathbf{\beta} + \mathbf{\varepsilon}$

- 2. Homoscedastic error variance: $Var(\varepsilon_i) = \sigma^2$.
- 3. Uncorrelated errors: $Cov(\varepsilon_i, \varepsilon_j) = 0$.
- 4. Additivity of errors: $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

We mainly use plots to assess this (more on model fit in F:3.4 Model choice and variable seletion)

- Covariate vs response (for each covariate)
- Covariate vs error (when we have simulated data and know the truth)
- Covariate vs residual (estimated error),
- Predicted response vs residual.

Uncorrelated errors?



Top: positively autocorrelated errors. Bottom: negatively correlated errors. Right: x vs y. Left: x vs error. Example from Fahrmeir et al (2013): Regression. Springer. (p.81). R code from TMA4267 lectures tab.



Fig. 3.4 Illustration for correlated residuals when the model is misspecified: Panel (**a**) displays (simulated) data based on the function $E(y_i | x_i) = \sin(x_i) + x_i$ and $\varepsilon_i \sim N(0, 0.3^2)$. Panel (**b**) shows the estimated regression line, i.e., the nonlinear relationship is ignored. The corresponding residuals can be found in panel (**c**)

Fahrmeir et al (2013): Regression. Springer. (p.82)

```
x1=runif(n,0,3)
x2=runif(n,0,3)
e=rnorm(n,0,0.4)
y=exp(1+x1-x2+e)
plot(x1,y,pch=20)
plot(x2,y,pch=20)
plot(x1,log(y),pch=20)
plot(x2,log(y),pch=20)
```

Multiplicative errors



Top: x1 and $x^{1/2}$ vs y. Bottom: x1 $a^{2/2}$ nd x2 vs log(y). Example from Fahrmeir et al (2013): Regression. Springer. (p.85). R code from TMA4267 lectures tab.

Covariates - how to include in the linear regression?

- 1. Continuous covariates: as is, transformed or using polynomials.
- 2. Categorical covariates: dummy variable or effect coding.
- 3. Interactions between covariates.
Munich rent index data

```
> colnames(ds)
[1] "rent" "rentsqm" "area" "yearc" "location" "bath"
[7] "kitchen" "cheating" "district"
> apply(ds[,1:4],2,summary)
          rent rentsqm area yearc
      40.51 0.4158 20.00 1918
Min.
1st Qu. 322.00 5.2610 51.00 1939
Median 427.00 6.9800 65.00 1959
Mean 459.40 7.1110 67.37 1956
3rd Qu. 559.40 8.8410 81.00 1972
Max.
       1843.00 17.7200 160.00 1997
> unlist(apply(ds[,5:8],2,table))
location.1 location.2 location.3 bath.0 bath.1 kitchen.0
     1794
                1210
                            78
                                    2891
                                           191
                                                 2951
kitchen.1 cheating.0 cheating.1
             321
                          2761
      131
```

How to code categorical covariates: rentsqm vs location with linear coding

Location average=1, good=2 and top=3, and regression model

rentsqm_i =
$$\beta_0 + \beta_1 \text{location}_i + \varepsilon_i$$

- Parameter estimate: $\hat{\beta}_1 = 0.39$. What does that mean?
 - Flat of average location: $\widehat{rentsqm} = \hat{\beta}_0 + \hat{\beta}_1 \cdot 1$
 - Flat of good location: rentsqm = \hat{\beta}_0 + \hat{\beta}_1 \cdot 2
 Flat of top location: rentsqm = \hat{\beta}_0 + \hat{\beta}_1 \cdot 3
- What is the difference in predicted rentsqm between top and good location, and between good and average location?
- So, the difference between a top and a good location is the same as the difference between good and average. Is this what we want?

Linear coding

Residual standard error: 2.427 on 3080 degrees of freedom Multiple R-squared: 0.007748,Adjusted R-squared: 0.007425 F-statistic: 24.05 on 1 and 3080 DF, p-value: 9.878e-07

rentsqm vs location with dummy variable coding

$$aloc_{i} = \begin{cases} 0 & location_{i} \text{ is not average} \\ 1 & location_{i} \text{ is average} \end{cases}$$
$$gloc_{i} = \begin{cases} 0 & location_{i} \text{ is not good} \\ 1 & location_{i} \text{ is good} \end{cases}$$
$$tloc_{i} = \begin{cases} 0 & location_{i} \text{ is not top} \\ 1 & location_{i} \text{ is top} \end{cases}$$

 $\mathsf{rentsqm}_i = \beta_0 + \beta_1 \mathsf{aloc}_i + \beta_2 \mathsf{gloc}_i + \beta_3 \mathsf{tloc}_i + \varepsilon_i$

- Write down the design matrix for this regression model, when we have 1794 flats with average location, 1210 with good and 78 with top location.
- What is the rank of this design matrix?
- Is there a problem, and a solution?

3.4 Dummy Coding for Categorical Covariates

For modeling the effect of a covariate $x \in \{1, ..., c\}$ with *c* categories using dummy coding, we define the c - 1 dummy variables

$$x_{i1} = \begin{cases} 1 & x_i = 1, \\ 0 & \text{otherwise,} \end{cases} \qquad \dots \qquad x_{i,c-1} = \begin{cases} 1 & x_i = c - 1, \\ 0 & \text{otherwise,} \end{cases}$$

for i = 1, ..., n, and include them as explanatory variables in the regression model

$$y_i = \beta_0 + \beta_1 x_{i1} + \ldots + \beta_{i,c-1} x_{i,c-1} + \ldots + \varepsilon_i.$$

For reasons of identifiability, we omit one of the dummy variables, in this case the dummy variable for category c. This category is called reference category. The estimated effects can be interpreted by direct comparison with the (omitted) reference category.

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.97)

Dummy coding via contr.treatment

```
> contrasts(ds$location)=contr.treatment(3)
```

> fit2=lm(rentsqm~location,data=ds)

```
> summary(fit2)
```

Call:

```
lm(formula = rentsqm ~ location, data = ds)
```

Coefficients:

	Estimate S	Std. Error	t value	Pr(> t)	
(Intercept)	6.95654	0.05728	121.456	< 2e-16	***
location2	0.31570	0.09025	3.498	0.000475	***
location3	1.21579	0.28060	4.333	1.52e-05	***
Signif. code	es: 0 '***	*' 0.001 '*	** 0.01	** 0.05	'.' 0.1 '

Residual standard error: 2.426 on 3079 degrees of freedom Multiple R-squared: 0.008867,Adjusted R-squared: 0.008223 F-statistic: 13.77 on 2 and 3079 DF, p-value: 1.109e-06

Effect coding via contr.sum

```
> contrasts(ds$location)=contr.sum(3)
> fit3=lm(rentsqm~location,data=ds)
> summary(fit3)
Call:
lm(formula = rentsqm ~ location, data = ds)
Coefficients:
           Estimate Std. Error t value Pr(>|t|)
(Intercept) 7.46704 0.09638 77.477 < 2e-16 ***
location1 -0.51050 0.10189 -5.010 5.75e-07 ***
location2 -0.19479 0.10445 -1.865 0.0623.
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Residual standard error: 2.426 on 3079 degrees of freedom Multiple R-squared: 0.008867,Adjusted R-squared: 0.008223 F-statistic: 13.77 on 2 and 3079 DF, p-value: 1.109e-06

Response: birth weight

Covariates: glucose level of mother and BMI of mother.



Figure from Kathrine Frey Frøslie.

Response: birth weight

Covariates: glucose level of mother and BMI of mother - with interaction.



Figure from Kathrine Frey Frøslie.

The classical linear model

$$\begin{array}{rcl} \mathbf{Y} & = & \mathbf{X} & \mathbf{\beta} & + & \mathbf{\varepsilon} \\ & & (n \times 1) \end{array} \\ E(\mathbf{\varepsilon}) &= & \mathbf{0} \\ & (n \times 1) \end{array} \quad \text{and} \quad Cov(\mathbf{\varepsilon}) &= & \sigma^2 \mathbf{I} \\ & & (n \times n) \end{array}$$

where

 \blacktriangleright β and σ^2 are unknown parameters and

• the design matrix **X** has *i*th row $[x_{i1}x_{i2}\cdots x_{ip}]$. Next: find the estimator $\hat{\beta}$.

- Model assessment: residual plots.
- Covariates: how to include in linear regression?
- Least squares and maximum likelihood estimator for β .

$$\hat{oldsymbol{eta}} = (oldsymbol{X}^{ op}oldsymbol{X})^{-1}oldsymbol{X}^{ op}oldsymbol{Y}$$

Covariates: how to include in the linear regression [F3.1.3]

Ex: Y = rentsorm
Xy = location

$$3: top$$

 $\beta_1 = 0.39$ \rightarrow the effect of good location
is twice the effect of overage location
a) Linear coding

$$X = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 3052x 4 & \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ \end{bmatrix} \begin{array}{c} \text{for average} \cdot 1774 \\ \text{good} \cdot 1210 \\ \text{top} \cdot 78 \end{array}$$

C) add a restriction: sum-to-zero

$$\begin{array}{c} 3\\ \overline{2}\\ \overline{2}\\ \overline{3}\\ \overline{2}\\ \overline{3}\\ \overline{5}\\ \overline{1}\\ \overline{5}\\ \overline{1}\\ \overline{5}\\ \overline{1}\\ \overline{5}\\ \overline{5}$$

$$Z_{12} \begin{cases} 1 & i \neq x_{1} = 1 \\ -1 & i \neq x_{1} = 3 \\ 0 & else(x_{1} = 2) \end{cases} \begin{cases} 1 & i \neq x_{1} = 2 \\ -1 & i \neq x_{1} = 3 \\ 0 & else(x_{1} = 2) \end{cases} \end{cases}$$

$$U_{1} = d_{0} + d_{1} Z_{1} + d_{2} \cdot Z_{2} + C$$

$$d_{3} = -d_{1} - d_{2}.$$

3) Interactions

- Is the effect (on Y) of a change in X1 dependent on the value of another covercede X2?
- Lego: Y = birth weight child X1: Glucose level of nother X2: BMI of nother (hog M2: BMI of nother (hog)

Figure 1: Y = po + py X1 + B2X2+ E

Figure 2: High glucose will have a different effect on birth weight when Brins low compared to when $Brin is high. <math>\Longrightarrow$ we have on interaction between X_1 and X_2 .

Eshmeter for
$$\beta$$
 [F3.2.1]
1) Maximum likelihood
If $\epsilon \sim N_n(0, \sigma^2 I)$ then $Y \sim N_n(X_p, \sigma^2 I)$
Alt 1: $Y_{1, Y_{2, \cdots}} Y_n$ independent
 $E(Y_i) = X_i^T \beta$, $Ver(Y_i) = \sigma^2$

$$\begin{aligned} x_{i}^{T}\beta \\ \rightarrow \int (y_{i}) y_{i}\sigma) &= \frac{1}{V_{2\pi}} \frac{1}{\sigma} \cdot e^{-\frac{1}{2}\sigma^{2}} (y_{i}-y_{i})^{2} \\ x_{i}^{T}\beta \\ \downarrow & x_{i}^{T}\beta \\ \end{pmatrix} \\ = \left(\frac{1}{2\pi}\right)^{2} \frac{1}{\sigma^{n}} \frac{1}{\varepsilon_{\pi}} \frac{1}{\sigma} \exp\left\{-\frac{1}{2\sigma^{2}} \frac{1}{\varepsilon_{\pi}} (y_{i}-x_{i}^{T}\beta)^{2}\right\} \\ + \log\left[\frac{1}{2\sigma}\right]^{2} \frac{1}{\sigma^{n}} \exp\left\{-\frac{1}{2\sigma^{2}} \frac{1}{\varepsilon_{\pi}} (y_{i}-x_{i}^{T}\beta)^{2}\right] \\ + \log\left[\frac{1}{2\sigma}\right]^{2} \frac{1}{\sigma^{n}} \exp\left\{-\frac{1}{2\sigma^{2}} \frac{1}{\varepsilon_{\pi}} (y_{i}-x_{i}^{T}\beta)^{2}\right\} \\ + \log\left[\frac{1}{2\sigma}\right]^{2} \frac{1}{\varepsilon_{\pi}} (y_{i}-x_{i}^{T}\beta)^{2} \frac{1}{\varepsilon_{\pi}} (y_{i}-x_{i}^{T}\beta)^{2}} \frac{1}{\varepsilon_{\pi}} (y_{i}-x_{i}^{T}\beta)^{2} \frac{1}{\varepsilon_{\pi}}$$

Alt 2:
$$Y \sim N_n(\mu, Z)$$

 $f(y; \mu, \Sigma) = (\frac{1}{2\pi})^{\frac{N}{2}} \left[\det(\Sigma) \right]^{-\frac{1}{2}}$
 $e_{xp} \left\{ -\frac{1}{2} (y-\mu)^T \Sigma^{-1} (y-\mu)^{2} \right\}$
Homework: $\mu = X_{F}, Z = 0^{e}T \Rightarrow get the Seme L(p, c^{e}) as (*)$

ii) To minimize
$$LS(p)$$
 with p we may solve
 $\frac{\partial LS(p)}{\partial p} = 0$
($\frac{\partial LS(p)}{\partial pTi}$)

"Need" two rules for don'velives
with vector:
 $\frac{\partial LS(p)}{\partial pTi}$

 $\frac{\partial LS(pTi}{\partial pTi}$

 $\frac{\partial LS(pTi}{\partial pTi}$

 $\frac{\partial LS(pTi}{\partial pTi}$

 $\frac{\partial LS(pTi}{\partial pTi}$

 $\frac{\partial LS($

Honework: Op LS(p) with these two rules!

X

TMA4267 Linear Statistical Models V2017 (L9) Part 2: Linear regression: Parameter estimation [F:3.2]

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: February 14, 2017

The classical linear model

$$\begin{array}{ll} \mathbf{Y} &=& \mathbf{X} \stackrel{\boldsymbol{\beta}}{}_{(n \times 1)} + \stackrel{\boldsymbol{\varepsilon}}{}_{(n \times 1)} \\ E(\boldsymbol{\varepsilon}) = \begin{array}{c} \mathbf{0} \\ (n \times 1) \end{array} \quad \text{and} \quad Cov(\boldsymbol{\varepsilon}) = \begin{array}{c} \sigma^2 \mathbf{I} \\ (n \times n) \end{array} \end{array}$$

where

 \blacktriangleright β and σ^2 are unknown parameters and

► the design matrix **X** has full rank, with *i*th row $[x_{i1}x_{i2}\cdots x_{ip}]$. Today

- 1. find estimator for β ,
- 2. find estimator for σ^2 , and
- 3. look at two idempotent matrices H and I H to arrive at
- 4. geometric interpretation.

Rules for derivatives with respect to a vector

- Let β be a *p*-dimensional column vector of interest,
- and let $\frac{\partial}{\partial\beta}$ denote the *p*-dimensional vector with partial derivatives wrt the *p* elements of β .
- Let *d* be a *p*-dimensional column vector of constants and
- **D** be a $p \times p$ symmetric matrix of constants.

Rule 1:

$$\frac{\partial}{\partial \boldsymbol{\beta}}(\boldsymbol{d}^{\mathsf{T}}\boldsymbol{\beta}) = \frac{\partial}{\partial \boldsymbol{\beta}}(\sum_{j=1}^{p} d_{j}\beta_{j}) = \boldsymbol{d}$$

Rule 2:

$$\frac{\partial}{\partial \boldsymbol{\beta}}(\boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{D}\boldsymbol{\beta}) = \frac{\partial}{\partial \boldsymbol{\beta}}(\sum_{j=1}^{p}\sum_{k=1}^{p}\beta_{j}d_{jk}\beta_{k}) = 2\boldsymbol{D}\boldsymbol{\beta}$$

See Härdle and Simes (2015), page 65, Equation (2.23) and (2.24).

Two questions

Have found least squares and maximum likelihood estimator for β :

$$\hat{\boldsymbol{eta}} = (\boldsymbol{X}^{\mathsf{T}} \boldsymbol{X})^{-1} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{Y}$$

and we have assumed that the rank(X) = p for $n \times p$ design matrix (where n > p).

- Q1: What can we say about $X^T X$?
- Q2: Why is the following wrong?

Using $(AB)^{-1} = B^{-1}A^{-1}$,

$$\hat{\boldsymbol{\beta}} = (\boldsymbol{X}^{T}\boldsymbol{X})^{-1}\boldsymbol{X}^{T}\boldsymbol{Y} = \boldsymbol{X}^{-1}(\boldsymbol{X}^{T})^{-1}\boldsymbol{X}^{T}\boldsymbol{Y} = \boldsymbol{X}^{-1}\boldsymbol{Y}$$

The classical linear model

The model

$$oldsymbol{Y} = oldsymbol{X}oldsymbol{eta} + arepsilon$$

is called a classical linear model if the following is true:

1.
$$E(\varepsilon) = 0$$
.

2.
$$\operatorname{Cov}(\varepsilon) = \operatorname{E}(\varepsilon \varepsilon^{T}) = \sigma^{2} I$$
.

3. The design matrix has full rank $rank(\mathbf{X}) = k + 1 = p$. The classical *normal* linear regression model is obtained if additionally

4. $\varepsilon \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$

holds. For random covariates these assumptions are to be understood conditionally on X.

Acid rain

occurs when emissions of sulfur dioxide (SO2) and oxides of nitrogen (NOx) react in the atmosphere with water, oxygen, and oxidants to form various acidic compounds. These compounds then fall to the earth in either dry form (such as gas and particles) or wet form (such as rain, snow, and fog).



Source: http://myecoproject.org/get-involved/pollution/acid-rain/



http://www.eoearth.org/view/article/149814/

Acid rain in Norwegian lakes

Measured pH in Norwegian lakes explained by content of

- ▶ x1: SO₄: sulfate (the salt of sulfuric acid),
- ► x2: N0₃: nitrate (the conjugate base of nitric acid),
- ► x3: *Ca*: calsium,
- ► x4: latent *AI*: aluminium,
- ► x5: organic substance,
- x6: area of lake,
- ► x7: position of lake (Telemark or Trøndelag),

pH is a measure of the acidity of alkalinity of water, expressed in terms of its concentration of hydrogen ions. The pH scale ranges from 0 to 14. A pH of 7 is considered to be neutral. Substances with pH of less that 7 are acidic; substances with pH greater than 7 are basic.



http://www.eoearth.org/view/article/149814/





0=Telemark, 1=Trondelag

Acid rain data



Output from fitting the full model in R

```
> fit=lm(y~.,data=ds)
> summary(fit)
Coefficients:
```

	E	lstima	te St	d.	Er	ror	t va	alue	e Pr	:(> t)			
(Intercep	ot) 5.	67643	34 ().13	389:	162	40	.862	2 <	< 2e-	16	**	*	
x1	-0.	31504	.44 ().0!	587	512	-5	. 362	24.	27e-	05	**	*	
x2	-0.	00185	33 (0.00	012	587	-1	.472	2	0.1	58			
xЗ	0.	97517	'45 ().14	449(075	6	.730) 2.	62e-	06	**	*	
x4	-0.	00022	.68 (0.00	010	038	-0.	.226	5	0.8	24			
x5	-0.	03342	.42 (0.02	225(009	-1	.485	5	0.1	55			
x6	-0.	00393	99 (0.0	724:	339	-0.	.054	ł	0.9	57			
x7	0.	08887	22 ().10	025	724	0	.866	5	0.3	98			
Signif. d	codes:	0 *	**' (0.00	01	,**,	0.0	D1 '	·*'	0.05	'.	,	0.1	,

Residual standard error: 0.1165 on 18 degrees of freedom Multiple R-squared: 0.93,Adjusted R-squared: 0.9027 F-statistic: 34.15 on 7 and 18 DF, p-value: 3.904e-09

Question: explain how to interpret $\hat{\beta}_0$ and $\hat{\beta}_3$.



3.10 Asymptotic Properties of the Least Squares Estimator

- 1. The least squares estimator $\hat{\beta}_n$ for β and the ML or REML estimator $\hat{\sigma}_n^2$ for the variance σ^2 are consistent.
- 2. The least squares estimator asymptotically follows a normal distribution, specifically

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}) \stackrel{d}{\rightarrow} \mathrm{N}(\boldsymbol{0}, \sigma^2 V^{-1}).$$

That is the difference $\hat{\beta}_n - \beta$ normalized with \sqrt{n} converges in distribution to the normal distribution on the right-hand side.

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.120)

Projection matrix: definition and properties

- A matrix **A** is a *projection matrix* if it is idempotent, $\mathbf{A}^2 = \mathbf{A}$.
- An idempotent matrix is an orthogonal projection matrix if, in the decomposition of a vector, v = Av + (v - Av), Av and v - Av = (I - A)v are always orthogonal, that is, (Av)^T(v - Av) = 0.
- A symmetric projection matrix is orthogonal.
- The eigenvalues of a projection matrix are 0 and 1.
- If a $(n \times n)$ symmetric projection matrix **A** has rank *r* then *r* eigenvalues are 1 and n r are 0.
- The trace and rank of a symmetric projection matrix are equal: tr(A) = rank(A).
Results so far

• Least squares and maximum likelihood estimator for β :

$$\hat{\boldsymbol{eta}} = (\boldsymbol{X}^{T} \boldsymbol{X})^{-1} \boldsymbol{X}^{T} \boldsymbol{Y}$$

• Restricted maximum likelihood estimator for σ^2 :

$$\hat{\sigma^2} = \frac{1}{n-p} (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})^T (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}}) = \frac{\mathsf{SSE}}{n-p}$$

Projection matrices: idempotent, symmetric/orthogonal:

$$\boldsymbol{H} = \boldsymbol{X} (\boldsymbol{X}^{T} \boldsymbol{X})^{-1} \boldsymbol{X}^{T}$$
$$\boldsymbol{I} - \boldsymbol{H} = \boldsymbol{I} - \boldsymbol{X} (\boldsymbol{X}^{T} \boldsymbol{X})^{-1} \boldsymbol{X}^{T}$$

with important connection:

$$\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$$

 $\hat{\mathbf{\varepsilon}} = \mathbf{I} - \mathbf{H}\mathbf{Y}$

Results from Mathematics 3

Best approximation theorem

The vector $\hat{\mathbf{Y}}$ in the column space of \mathbf{X} that makes $|| \mathbf{Y} - \hat{\mathbf{Y}} ||$ as small as possible, is the orthogonal projection of \mathbf{Y} on the column space of \mathbf{X} .

Orthogonal decomposition

We want $\hat{\boldsymbol{\beta}}$ to minimize $|| \boldsymbol{Y} - \hat{\boldsymbol{Y}} || = (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})^T (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})$ (least squares principle).

The column space of X consists of vectors of the form $X\hat{\beta}$, so $X\hat{\beta}$ is the orthogonal projection of Y onto the column space of X. $\hat{Y} = HY$, and $H = X(X^TX)^{-1}X^T$ projects onto the column space of X. Observe: HX = X.

This is equivalent to observing that $\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ is in the orthogonal complement of the column space of \mathbf{X} .

 $\hat{\varepsilon} = \mathbf{Y} - \mathbf{H}\mathbf{Y} = (\mathbf{I} - \mathbf{H})\mathbf{Y}$, and $\mathbf{I} - \mathbf{H}$ projects onto the space orthogonal to the column space of \mathbf{X} . Observe: $(\mathbf{I} - \mathbf{H})\mathbf{X} = \mathbf{0}$

That is, $\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ is orthogonal to all columns of \mathbf{X} , so $\mathbf{X}^{T}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{0}$ and $\mathbf{X}^{T}\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}^{T}\mathbf{Y}$.



Putanen, Styan and Isotalo: Matrix Tricks for Linear Statistical Models: Our Personal Top Twenty, Figure 8.3.

3.7 Geometric Properties of the Least Squares Estimator

The method of least squares has the following geometric properties:

- 1. The predicted values \hat{y} are orthogonal to the residuals $\hat{\varepsilon}$, i.e., $\hat{y}'\hat{\varepsilon} = 0$.
- 2. The columns x^{j} of X are orthogonal to the residuals $\hat{\boldsymbol{\varepsilon}}$, i.e., $(x^{j})'\hat{\boldsymbol{\varepsilon}} = 0$ or $X'\hat{\boldsymbol{\varepsilon}} = \mathbf{0}$.
- 3. The average of the residuals is zero, i.e.,

$$\sum_{i=1}^{n} \hat{\varepsilon}_i = 0 \quad \text{or} \quad \frac{1}{n} \sum_{i=1}^{n} \hat{\varepsilon}_i = 0.$$

4. The average of the predicted values \hat{y}_i is equal to the average of the observed response y_i , i.e.,

$$\frac{1}{n}\sum_{i=1}^n \hat{y}_i = \bar{y}.$$

5. The regression hyperplane runs through the average of the data, i.e.,

$$\bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}_1 + \dots + \hat{\beta}_k \bar{x}_k.$$

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.112) Alternative summery of Geometry of Least Squares

- Mean response vector: $E(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta}$
- As *β* varies, *Xβ* spans the model plane of all linear combinations. I.e. the space spanned by the columns of *X*: the column-space of *X*.
- Due to random error (and unobserved covariates), Y is not exactly a linear combination of the columns of X.
- LS-estimation chooses $\hat{\beta}$ such that $X\hat{\beta}$ is the point in the column-space of X that is closes to Y.
- The residual vector $\hat{\boldsymbol{\varepsilon}} = \boldsymbol{Y} \hat{\boldsymbol{Y}} = (\boldsymbol{I} \boldsymbol{H})\boldsymbol{Y}$ is perpendicular to the column-space of \boldsymbol{X} .
- Multiplication by H = X(X^TX)⁻¹X^T projects a vector onto the column-space of X.
- Multiplication by I H = I X(X^TX)⁻¹X^T projects a vector onto the space perpendicular to the column-space of X.

Today

• Least squares and maximum likelihood estimator for β :

$$\hat{oldsymbol{eta}} = (oldsymbol{X}^{ op}oldsymbol{X})^{-1}oldsymbol{X}^{ op}oldsymbol{Y}$$

has mean $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$ and $Cov(\hat{\boldsymbol{\beta}}) = \sigma^2 (\boldsymbol{X}^T \boldsymbol{X})^{-1}$.

- For the normal model: $\hat{\boldsymbol{\beta}} \sim N_p(\boldsymbol{\beta}, \sigma^2(\boldsymbol{X}^T \boldsymbol{X})^{-1})$.
- Asymptotic properties of the least squares estimator: normality.
- Orthogonal projection matrices *H* and *I H* with geometric interpretation.

Next time: properties of residuals and $\hat{\sigma}^2$, confidence intervals and hypothesis testing for regression coefficients.

$$Y = X\beta + \xi , E(\xi) = 0, Ca(\xi) = 0^{2} I$$

$$Inter Inter Inte$$

iii) solving the normal equations (substitutes with
$$\beta$$
)

$$\frac{\hat{\beta}}{\hat{\beta}} = (X^T X)^{-1} X^T Y$$

Questions on slide: X renup: XTX pxp motivx symmetric positive definite uverse exists

If UT XTXU > 0 for all UFO then XTX posibive definite.

 $0 \leq (vX)^{\gamma}(vX) = vXX^{\gamma}v$

Assume that UT XTXV = 0 than XV=0. If X has full reach than XV=0 only has V=0 as solution. => then XTX must be possive definite.

\$=(XTX)-187 is the less' squared enhanctor of p. If we assume a normal linear model then \$ 15 also the maximum likelihood enhanctor Ex: Acid roin: fit in R wing Im

St. Error: SD(B)

4

Properties of
$$\beta$$

 $\beta = (XTX)^{-1}XTY$ and $E(Y) = Xp$
 C RV
 $Gar(Y) = \sigma^{2} I$
 $Gar(Y) = \sigma^{2} I$

Find
$$E(\beta)$$
 and $Gv(\beta)$:
 $E(\beta) = E(CY) = CE(Y) = (XTX)^T X = B$
 $X = T$
 $X = T$
 $X = T$
 $X = T$
 $X = T$

In a normal model: pha Np(p, 02 (XTX)-1)

EX: Acid rain: Whet is
$$Ve-(\beta_3)?$$

 $Sp(p_3) = 0.144: (\hat{\sigma}^2 \cdot (XTX)'_{24,42})$

5

Last information on
$$\beta$$
: From part 1:
 $(\beta - E(\beta))^{T} \text{GV}(\beta)^{T}(\beta - E(\beta)) \sim X_{p}$
 $\overline{d^{2}}(\beta - \beta)^{T}(X^{T}X)(\beta - \beta) \sim X_{p}$
 $G_{v}(\beta)^{2} = 0^{v}(X^{T}X)^{1}$
 $G_{v}(\beta)^{1} = \overline{d^{2}}(X^{T}X)$
 $\overline{G_{v}(\beta)^{1}} = \overline{d^{2}}(X^{T}X)$
 $\overline{G_{v}(\beta)^{1}} = \overline{d^{2}}(X^{T}X)$
 $\overline{G_{v}(\beta)^{1}} = \overline{d^{2}}(X^{T}X)$
 $\overline{G_{v}(\beta)^{1}} = (\overline{d^{2}})^{\frac{N}{2}} (\overline{d^{2}})^{\frac{N}{2}} \exp\{-\overline{2\sigma^{2}}(y - X_{p})^{T}(y - X_{p})\}$
 $L(\beta, \sigma^{2}) = \ln(L(\beta, \sigma^{2}))$
 $= -\frac{n}{2}\ln(RT) - \frac{n}{2}\ln\sigma^{2} - \frac{1}{2\sigma^{2}}(y - X_{p}^{2})^{T}(y - X_{p}^{2})$

Need estimater for (T.)

$$\frac{\partial l}{\partial \sigma^2} = 0 \quad (=) \qquad \qquad \frac{\partial l (n \times \frac{1}{2})}{x} = \frac{\partial l \times \frac{1}{2}}{y} = \frac{\partial l \times \frac{1}{2}}{y} = 0$$

$$\frac{\partial l (n \times \frac{1}{2})}{x} = \frac{\partial l \times \frac{1}{2}}{y} = 0$$

6

$$\frac{\Pi}{\sigma^2} = \frac{1}{\sigma^4} \left(y - X_{j}^2 \right)^{\dagger} \left(y X_{j}^2 \right)$$

$$\frac{\int_{\sigma_{1L}}^2 - \frac{1}{\sigma^4} \left(Y - X_{j}^2 \right)^{\dagger} \left(Y - X_{j}^2 \right)}{\hat{\epsilon}} = n \hat{\epsilon}^{\dagger} \hat{\epsilon}^{\dagger} \hat{\epsilon}^{\dagger}$$

$$\hat{\epsilon}^{\dagger} \hat{\epsilon} = 3ums of squees of errors SSE$$

But, this exhiber is rerely used, because it is brased (use tr-formula Pert 1 to jund the meen) However: is unbased

towever:

$$\hat{\sigma}^2 = \frac{1}{n-p} \hat{\varepsilon}^{+} \hat{\varepsilon}^{-}$$
is unbased

THA 1295 Stetistical Inforence ; more on this.

Ex: Acid rain
$$\widehat{SD}(\widehat{\beta}_{2}) = \sqrt{(272)^{\frac{1}{2}} + \sqrt{3} \cdot \widehat{S}^{2}}$$

 $\widehat{S}^{2} = n - p \widehat{E}^{+} \widehat{E}$ $\widehat{S}^{2} = 0.1165$
 \widehat{S}^{1} is residual stendard error in printout

Z

$$E(Y) = X_{fS}, so E(Y) = X_{fS}^{2} = Y \leftarrow prediction$$

$$fS = (X^{T}X)^{-1}X^{T}Y$$

$$fS = X_{fS}^{2} = X(X^{T}X)^{-1}X^{T}Y = HY$$

$$H$$

H= $X(X^TX)^{-1}X^T$ is called the "hat matrix" nown for putting the last on Y.

Observe (see also RecEr3. P3a) that H is symmetric Is idemposent H² = H has renh p = show this (I-H) is also symmetric and idemposent, and renh (I-H) = n-p. = show this

8

TMA4267 Linear Statistical Models V2017 (L10) Part 2: Linear regression: Parameter estimation [F:3.2], Properties of residuals and distribution of estimator for error variance Confidence interval and hypothesis for one regression coefficient

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: February 17, 2017

- 1. Properties for residuals (from the hat matrix), leading to properties for $\hat{\sigma}^2$,
- 2. Then, confidence interval and hypothesis test for regression coefficient.

The classical linear model

The model

$$\mathbf{Y} = \mathbf{X}oldsymbol{eta} + arepsilon$$

is called a classical linear model if the following is true:

1.
$$E(\varepsilon) = 0$$
.

2.
$$\operatorname{Cov}(\varepsilon) = \operatorname{E}(\varepsilon \varepsilon^{T}) = \sigma^{2} I$$
.

3. The design matrix has full rank $rank(\mathbf{X}) = k + 1 = p$. The classical *normal* linear regression model is obtained if additionally

1. $\varepsilon \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$

holds. For random covariates these assumptions are to be understood conditionally on X.

Results so far

• Least squares and maximum likelihood estimator for β :

$$\hat{oldsymbol{eta}} = (oldsymbol{X}^{ op}oldsymbol{X})^{-1}oldsymbol{X}^{ op}oldsymbol{Y}$$

with mean $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$ and $Cov(\hat{\boldsymbol{\beta}}) = \sigma^2 (\boldsymbol{X}^T \boldsymbol{X})^{-1}$.

• Restricted maximum likelihood estimator for σ^2 :

$$\hat{\sigma^2} = \frac{1}{n-p} (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})^T (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}}) = \frac{SSE}{n-p}$$

Projection matrices: idempotent, symmetric/orthogonal:

$$\boldsymbol{H} = \boldsymbol{X}(\boldsymbol{X}^{T}\boldsymbol{X})^{-1}\boldsymbol{X}^{T}$$

projects onto column space of \boldsymbol{X}

$$\boldsymbol{I} - \boldsymbol{H} = \boldsymbol{I} - \boldsymbol{X} (\boldsymbol{X}^{\mathsf{T}} \boldsymbol{X})^{-1} \boldsymbol{X}^{\mathsf{T}}$$

projects onto space orthogonal to column space of \boldsymbol{X}

with important connection: predictions $\hat{Y} = HY$ and residuals $\hat{\varepsilon} = (I - H)Y$



Putanen, Styan and Isotalo: Matrix Tricks for Linear Statistical Models: Our Personal Top Twenty, Figure 8.3.

Quadratic forms [F:B3.3, Theorem B.2]

Random vector \boldsymbol{X} with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$, symmetric constant matrix \boldsymbol{A} .

• Quadratic form: $X^T A X$.

► The "trace-formula": $E(X^T A X) = tr(A \Sigma) + \mu^T A \mu$. Then, let $X \sim N_p(0, I)$, and R is a symmetric and idempotent matrix with rank r.

$$oldsymbol{X}^{T}oldsymbol{R}oldsymbol{X}\sim\chi^2_r$$

Now, also **S** is a symmetric and idempotent matrix with rank *s*, and RS = 0.

$$\frac{s \boldsymbol{X}^T \boldsymbol{R} \boldsymbol{X}}{r \boldsymbol{X}^T \boldsymbol{S} \boldsymbol{X}} \sim F_{r,s}$$

Properties: $\hat{\boldsymbol{\beta}}$ and $\hat{\sigma}^2$

• Least squares and maximum likelihood estimator for β :

$$\hat{oldsymbol{eta}} = (oldsymbol{X}^{ op}oldsymbol{X})^{-1}oldsymbol{X}^{ op}oldsymbol{Y}$$

has mean $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$ and $Cov(\hat{\boldsymbol{\beta}}) = \sigma^2 (\boldsymbol{X}^T \boldsymbol{X})^{-1}$.

- In addition is best linear unbiased estimator (BLUE), that is, among all unbiased estimator it has minimum variance in each component. (More in TMA4295 Statistical Inference.)
- For the normal model: $\hat{\boldsymbol{\beta}} \sim N_p(\boldsymbol{\beta}, \sigma^2(\boldsymbol{X}^T \boldsymbol{X})^{-1}).$
- Restricted maximum likelihood estimator for σ^2 :

$$\hat{\sigma^2} = \frac{1}{n-p} (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})^T (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}}) = \frac{\mathsf{SSE}}{n-p}$$

For the normal model

$$\frac{(n-p)\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-p}$$

Acid rain in Norwegian lakes

Measured pH in Norwegian lakes explained by content of

- ▶ x1: SO₄: sulfate (the salt of sulfuric acid),
- ► x2: N0₃: nitrate (the conjugate base of nitric acid),
- ► x3: *Ca*: calsium,
- ► x4: latent AI: aluminium,
- ► x5: organic substance,
- ► x6: area of lake,
- ► x7: position of lake (Telemark or Trøndelag),

Random sample of n = 26 lakes.

Output from fitting the full model in R

```
> fit=lm(y~.,data=ds)
> summary(fit)
```

Coefficients:

	E	sti	mate	Std.	Erro	or t	valı	ıe P	r(> t)		
(Intercep	t) 5.	676	4334	0.1	38916	52	40.86	52	< 2e-3	16	**	*
x1	-0.	315	0444	0.0	58751	.2	-5.36	52 4	.27e-0	05	**	*
x2	-0.	001	8533	0.0	01258	87	-1.47	72	0.1	58		
xЗ	0.9	975	1745	0.1	44907	'5	6.73	30 2	.62e-0	26	**	*
x4	-0.	000	2268	0.0	01003	88	-0.22	26	0.8	24		
x5	-0.	033	4242	0.0	22500	9	-1.48	35	0.1	55		
x6	-0.	003	9399	0.0	72433	89	-0.05	54	0.9	57		
x7	0.0	088	8722	0.1	02572	24	0.86	56	0.39	98		
Signif. c	odes:	0	·***	0.0	01 '*	**'	0.01	,*,	0.05	'.	,	0.1

Residual standard error: 0.1165 on 18 degrees of freedom Multiple R-squared: 0.93,Adjusted R-squared: 0.9027 F-statistic: 34.15 on 7 and 18 DF, p-value: 3.904e-09

, , 1

W. S. Gosset alias Student



Historical: Student-t fordelingen

- W.S. Gosset (1876-1937) was employed by the Guinness Brewing Company of Dublin.
- Sample sizes available for experimentation in brewing were necessarily small, and Gosset knew that a correct way of dealing with small samples was needed.
- He consulted Karl Pearson (1857-1936) of University College in London about the problem. Pearson told him the current state of knowledge was unsatisfactory.
- The following year Gosset undertook a course of study under Pearson. An outcome of his study was the publication in 1908 of Gosset's paper on "The Probable Error of a Mean," which introduced a form of what later became known as Student's t-distribution.
- Gosset's paper was published under the pseudonym "Student."
- The modern form of Student's t-distribution was derived by R.A. Fisher and first published in 1925.

t-distribution



DEF: *t*-distribution

Let Z be a standard normal random variable and V a chi-squared random variable with parameter ν (degrees of freedom). If Z and V are independent, the distribution of the random variable T

$$T = \frac{Z}{\sqrt{V/\nu}}$$

has probability density function

$$h(t) = \frac{\Gamma[(\nu+1)/2]}{\Gamma(\nu/2)\sqrt{\pi\nu}} (1 + \frac{t^2}{\nu})^{-(\nu+1)/2}$$

for $-\infty < t < \infty$. This distribution is called the (Student) *t*-distribution with ν degrees of freedom.

•
$$E(T) = 0$$
 if $\nu \geq 2$.

•
$$\operatorname{Var}(T) = \frac{\nu}{\nu - 2}$$
 if $\nu \geq 3$.

Are $\hat{\beta}$ and SSE are independent?

Independence – from Part 1: Let $X_{(p \times 1)}$ be a random vector from $N_p(\mu, \Sigma)$. Then AX and BX are independent iff $A\Sigma B^T = 0$.

We have:

$$\blacktriangleright \boldsymbol{Y} \sim N_n(\boldsymbol{X}\boldsymbol{\beta},\sigma^2\boldsymbol{I})$$

•
$$\boldsymbol{A} \boldsymbol{Y} = \hat{\boldsymbol{\beta}} = (\boldsymbol{X}^{T} \boldsymbol{X})^{-1} \boldsymbol{X}^{T} \boldsymbol{Y}$$
, and

$$\blacktriangleright BY = (I - H)Y.$$

Now
$$\boldsymbol{A}\sigma^2\boldsymbol{I}\boldsymbol{B}^T = \sigma^2\boldsymbol{A}\boldsymbol{B}^T = \sigma^2(\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T(\boldsymbol{I}-\boldsymbol{H}) = \boldsymbol{0}$$

- ► since X(I H) = X HX = X X = 0.
- We conclude that $\hat{oldsymbol{eta}}$ is independent of (I H)Y,
- ▶ and, since SSE=function of (I H)Y: SSE= $Y^T(I H)Y$,
- then $\hat{\beta}$ and SSE are independent.

Quantiles and critical values: N og t: $\alpha/2 = 0.025$



Kritiske verdier i *t*-fordelingen

 $P(T > t_{\alpha,\nu}) = \alpha$

ν	$' \ \alpha$.150	.100	.075	.050	.025	.010	.005	.001	.0005
	1	1.963	3.078	4.165	6.314	12.706	31.821	63.657	318.309	636.619
	2	1.386	1.886	2.282	2.920	4.303	6.965	9.925	22.327	31.599
	3	1.250	1.638	1.924	2.353	3.182	4.541	5.841	10.215	12.924
	4	1.190	1.533	1.778	2.132	2.776	3.747	4.604	7.173	8.610
	5	1.156	1.476	1.699	2.015	2.571	3.365	4.032	5.893	6.869
	6	1.134	1.440	1.650	1.943	2.447	3.143	3.707	5.208	5.959
	7	1.119	1.415	1.617	1.895	2.365	2.998	3.499	4.785	5.408
	8	1.108	1.397	1.592	1.860	2.306	2.896	3.355	4.501	5.041
	9	1.100	1.383	1.574	1.833	2.262	2.821	3.250	4.297	4.781
	10	1.093	1.372	1.559	1.812	2.228	2.764	3.169	4.144	4.587
	11	1.088	1.363	1.548	1.796	2.201	2.718	3.106	4.025	4.437
	12	1.083	1.356	1.538	1.782	2.179	2.681	3.055	3.930	4.318
	13	1.079	1.350	1.530	1.771	2.160	2.650	3.012	3.852	4.221
	14	1.076	1.345	1.523	1.761	2.145	2.624	2.977	3.787	4.140
	15	1.074	1.341	1.517	1.753	2.131	2.602	2.947	3.733	4.073
	16	1.071	1.337	1.512	1.746	2.120	2.583	2.921	3.686	4.015
	17	1.069	1.333	1.508	1.740	2.110	2.567	2.898	3.646	3.965
	18	1.067	1.330	1.504	1.734	2.101	2.552	2.878	3.610	3.922
	19	1.066	1.328	1.500	1.729	2.093	2.539	2.861	3.579	3.883
	20	1.064	1.325	1.497	1.725	2.086	2.528	2.845	3.552	3.850
	21	1.063	1.323	1.494	1.721	2.080	2.518	2.831	3.527	3.819
	22	1.061	1.321	1.492	1.717	2.074	2.508	2.819	3.505	3.792
	23	1.060	1.319	1.489	1.714	2.069	2.500	2.807	3.485	3.768
	24	1.059	1.318	1.487	1.711	2.064	2.492	2.797	3.467	3.745
	25	1.058	1.316	1.485	1.708	2.060	2.485	2.787	3.450	3.725
	26	1.058	1.315	1.483	1.706	2.056	2.479	2.779	3.435	3.707
	27	1.057	1.314	1.482	1.703	2.052	2.473	2.771	3.421	3.690
	28	1.056	1.313	1.480	1.701	2.048	2.467	2.763	3.408	3.674
	29	1.055	1.311	1.479	1.699	2.045	2.462	2.756	3.396	3.659
	30	1.055	1.310	1.477	1.697	2.042	2.457	2.750	3.385	3.646
	35	1.052	1.306	1.472	1.690	2.030	2.438	2.724	3.340	3.591
	40	1.050	1.303	1.468	1.684	2.021	2.423	2.704	3.307	3.551
	50	1.047	1.299	1.462	1.676	2.009	2.403	2.678	3.261	3.496
	60	1.045	1.296	1.458	1.671	2.000	2.390	2.660	3.232	3.460
	80	1.043	1.292	1.453	1.664	1.990	2.374	2.639	3.195	3.416
1	00	1.042	1.290	1.451	1.660	1.984	2.364	2.626	3.174	3.390
1	120	1.041	1.289	1.449	1.658	1.980	2.358	2.617	3.160	3.373
	∞	1.036	1.282	1.440	1.645	1.960	2.326	2.576	3.090	3.291

Acid rain in R

ds=read.table("https://www.math.ntnu.no/emner/ TMA4267/2017v/acidrain.txt",header=TRUE) fit=lm(y~.,data=ds) > confint(fit)

	2.5 %	97.5 %
(Intercept)	5.384581378	5.9682854281
x1	-0.438476153	-0.1916126966
x2	-0.004497716	0.0007911594
x3	0.670735075	1.2796138706
x4	-0.002335625	0.0018820903
x5	-0.080696921	0.0138484550
x6	-0.156117992	0.1482381575
x7	-0.126624544	0.3043688780

P-values: http://www.statistrikk.no/wp-content/uploads/ 2017/02/nerdekort.jpg

Today

- Distribution of SSE/ σ^2 is chisquared (n p).
- Independence of $\hat{\beta}$ and SSE.
- Inference about *β* components can be performed using the *t*-distribution

Distribution of SSE and
$$\hat{\sigma}^{2}$$

SSE = $\hat{\varepsilon}T\hat{\varepsilon} = Y^{T}(T-H)(T-H)Y$
 $\hat{\tau}$ $(y_{i} - \hat{y})^{\varepsilon}$ $\hat{\varepsilon} = (T-H)Y$
 $\hat{\varepsilon}$ $\hat{\varepsilon}$ $(T-H)Y$
remainder Fenn $(T-H) = n-p$.
Rec $(\xi, 3, P3)$ looks at the distribution of $\hat{\sigma}^{2} YT(T-H)Y = \hat{\sigma}^{2}$
by using the result on questric forms from Pert 1
(see stick)
 $Y \sim N_{n}(X_{i}\beta, \sigma^{\varepsilon}T)$
 $Y^{*} = \hat{\sigma} ((T-X_{i}\beta) \sim N_{n}(O, T))$
 $Y^{*T}(T-H)Y^{*} = \hat{\sigma}^{2} Y^{T}(T-H)Y$
 $(T-H)(Y-\hat{\Sigma}p)$ H
 $(T-H)Y - (T-H)X_{i}\beta$ $\hat{S} = \hat{\sigma}^{2} - (X_{n-p}^{2})$
 $\hat{V} = \hat{S} = \hat{\sigma}^{2} = \hat{\sigma}^{2} - \hat{X}_{n-p}^{2}$
 $\hat{V} = \hat{S} = \hat{\sigma}^{2} = \hat{\sigma}^{2} - \hat{X}_{n-p}^{2}$

$$E(U) = n-p \quad Vor(V) = 2(n-p)$$

$$Is \quad \hat{\sigma}^{2} \text{ an unbiased estimator }^{2}$$

$$E(\hat{\sigma}^{2}) = E(\frac{1}{n-p} \text{ SSE}) = n-p \quad E(\sigma^{2}V)$$

$$= \frac{\sigma^{2}}{n-p} \quad E(V) = \sigma^{2} \quad \text{unbiased.}$$
This is the when we assume $c = N$.

If we do not assume
$$E \cap N$$
, then we can
use the trece-formula
 $E(SSE) = E(YT(I-H)Y)$ $E(Y) = X_{fS}$
 $Ga(Y) = \sigma^{2}T$
 $= fr((T-H) \sigma^{2}T) + (X_{fS})^{T}(I+H) X_{fS}$
 $= (n-p)\sigma^{2} + 0$
 $E(\Im^{2}) - E(\frac{SSE}{n-p}) = \sigma^{2}$

Inference ebout one Bj

Ex: Acid (an
$$B_1 = effect of SOy on ptt of lake
$$\begin{array}{l}
\beta_1 = -0.315 \\
SD(\beta_1) = \sqrt{5^2 \left[(2T2)^{-1} \right]} \left[diagonalden \right] \\
SD(\beta_1) = \sqrt{5^2 \left[(2T2)^{-1} t_{X_1 X_1} \right]^T} \quad corresp. to Soy \\
X_1 \\
\end{bmatrix} \\
St. Error \Rightarrow 0.0587 in printat$$

$$\begin{array}{l}
\delta : "Residual chenderderror" = 0.1165 \\
\sqrt{556^T} \\
n=26 \\
n=26 \\
n=8 \\
\end{array} n-p = 18$$
"on 18 degrees of breedom".$$

To find a confidence inhual for Bj -or to test hypotheses about B; we need to know the distribution of a stabilic involving fij and Bjwith no other kinknown parameteo.

$$\hat{\beta}_{j} \sim N_{1}(\beta_{j}, G^{2}(X^{T}X)^{-1}_{U_{j}})$$

4
and
$$\int_{2}^{2} = \frac{SSE}{n-p}$$
 when $\frac{SSE}{\sigma^{2}} \sim N^{2}_{n-p}$.
Then:
 $\frac{\beta_{j} - \beta_{j}}{\sqrt{S_{j}\sigma^{2}}} \sim N(0,1)$
 $\sqrt{S_{j}\sigma^{2}} \sim SO(\beta_{j})$
 $\int_{2}^{2} \int_{2}^{2} \int_{2}^{$

$$T_j = \frac{\beta_j - \beta_j}{\sqrt{c_{jj}} \hat{c}} \wedge t_{n-p}$$

General result:

$$\frac{N(0,1)}{\sqrt{\frac{\chi^2_q}{q}}} \sim t_q$$

Ś

We have: $\frac{\beta_j - \beta_j}{\sqrt{c_{jj}}} \sim N(0,1)$ $\sqrt{c_{jj}} = \sigma$ and $(n-p)\frac{\delta^2}{\sigma^2} \sim \chi^2 n-p$



$$\hat{\beta}_{j}$$
 end $\hat{\sigma}^{2}$ need to be independent incrdue
that this holds. $\Rightarrow \text{Rec} \in 3.P3 + \text{slides}$
Use $T_{j} = \frac{\hat{\beta}_{j} - \hat{\beta}_{j}}{\sqrt{E_{j}}} \sim t_{n-p}$ for inference.
(C) Find a 25% confidence interval (CCI)
for $\hat{\beta}_{j}$

b



$$E_{X}: \begin{array}{c} A \\ B_{1} = -0.815 \\ \sqrt{c_{jj}} C_{j} = 0.058 \\ n=26, r=8 \\ f_{0.025, 18} = 2.1 \end{array} \qquad -0.815 \pm 2.1 \cdot 0.058 \\ = \underline{\Gamma - 0.44}, -0.19 \\ \end{array}$$

We see that 0 is not in the interval - what does this meen? \Rightarrow Reject Ho: $g_{j}=0$ is Hi $g_{j}\neq 0$ at sign level 5%

TMA4267 Linear Statistical Models V2017 (L11) Part 2: Linear regression:

Parameter estimation [F:3.2] and model selection [F:3.4] Hypothesis test for one regression coefficient Studentized and standardized residuals decomposition of variability and significance of regression R^2 , SPSE=Expected squared prediction error

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: February 21, 2017

Today

- 1. Hypothesis testing for β_j .
- 2. Residuals: standardized (or studentized) preferred.
- 3. Decomposition of variability: SST=SSR+SSE, and significance of regression.
- 4. R^2 gives the proportion of variability explained by the regression model. and will never decrease if new covariates are added to the model.
- 5. Model choice considerations.
- 6. SPSE: Expected squared prediction error.

The classical linear model

The model

$$\mathbf{Y} = \mathbf{X}oldsymbol{eta} + arepsilon$$

is called a classical linear model if the following is true:

1.
$$E(\varepsilon) = 0$$
.

2.
$$\operatorname{Cov}(\varepsilon) = \operatorname{E}(\varepsilon \varepsilon^{T}) = \sigma^{2} I$$
.

3. The design matrix has full rank $rank(\mathbf{X}) = k + 1 = p$. The classical *normal* linear regression model is obtained if additionally

1. $\varepsilon \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$

holds. For random covariates these assumptions are to be understood conditionally on X.

Properties for the normal linear model

• Least squares and maximum likelihood estimator for β :

$$\hat{oldsymbol{eta}} = (oldsymbol{X}^{ op}oldsymbol{X})^{-1}oldsymbol{X}^{ op}oldsymbol{Y}$$

with $\hat{\boldsymbol{\beta}} \sim N_{p}(\boldsymbol{\beta}, \sigma^{2}(\boldsymbol{X}^{T}\boldsymbol{X})^{-1}).$

• Restricted maximum likelihood estimator for σ^2 :

$$\hat{\sigma^2} = \frac{1}{n-p} (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}})^T (\boldsymbol{Y} - \boldsymbol{X}\hat{\boldsymbol{\beta}}) = \frac{SSE}{n-p}$$

with $\frac{(n-p)\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-p}$.

Statistic for inference about β_j, c_{jj} is diagonal element j of (X^TX)⁻¹.

$$T_j = \frac{\hat{\beta}_j - \beta_j}{\sqrt{c_{jj}}\hat{\sigma}} \sim t_{n-p}$$

Acid rain in Norwegian lakes

Measured pH in Norwegian lakes explained by content of

- ▶ x1: SO₄: sulfate (the salt of sulfuric acid),
- ► x2: N0₃: nitrate (the conjugate base of nitric acid),
- ► x3: *Ca*: calsium,
- ► x4: latent AI: aluminium,
- ► x5: organic substance,
- ► x6: area of lake,
- ► x7: position of lake (Telemark or Trøndelag),

Random sample of n = 26 lakes.

Output from fitting the full model in R

```
> fit=lm(y~.,data=ds)
> summary(fit)
```

Coefficients:

	E	sti	mate	Std.	Erro	or t	valı	ıe P	r(> t)		
(Intercep	t) 5.0	676	4334	0.1	38916	52	40.86	52	< 2e-3	16	**	*
x1	-0.	315	0444	0.0	58751	.2	-5.36	52 4	.27e-0	05	**	*
x2	-0.	001	8533	0.0	01258	87	-1.47	72	0.1	58		
xЗ	0.9	975	1745	0.1	44907	'5	6.73	30 2	.62e-0	26	**	*
x4	-0.	000	2268	0.0	01003	88	-0.22	26	0.8	24		
x5	-0.	033	4242	0.0	22500	9	-1.48	35	0.1	55		
x6	-0.	003	9399	0.0	72433	89	-0.05	54	0.9	57		
x7	0.0	088	8722	0.1	02572	24	0.86	56	0.39	98		
Signif. c	odes:	0	·***	0.0	01 '*	**'	0.01	,*,	0.05	'.	,	0.1

Residual standard error: 0.1165 on 18 degrees of freedom Multiple R-squared: 0.93,Adjusted R-squared: 0.9027 F-statistic: 34.15 on 7 and 18 DF, p-value: 3.904e-09

, , 1

Quantiles and critical values: N og t: $\alpha/2 = 0.025$



In R: specify area to the left, but our notation gives area to the right. Fahrmeir et al: notation with area to the left.

Properties of the residuals

- Residuals (raw): $\hat{\boldsymbol{\varepsilon}} = \boldsymbol{Y} \hat{\boldsymbol{Y}}$.
- with mean $E(\hat{\varepsilon}) = 0$ and covariance matrix $Cov(\hat{\varepsilon}) = \sigma^2(I - H)$ where $H = X(X^T X)^{-1} X^T$.
- In the normal model ε ~ N_n(0, σ²I) and then also the vector of residuals are normal, but with heteroscedastic variances and non-zero covariances.
- Standardized residuals: divide (raw) residuals by estimated standard deviation.
- Studentized residuals: leave-one-out version.
- Studentized residuals are compared with the normal distribution to assess normality of the error term.

3.12 Overview of Residuals

Ordinary Residuals

The residuals are given by

$$\hat{\varepsilon}_i = y_i - \hat{y}_i = y_i - \boldsymbol{x}'_i \hat{\boldsymbol{\beta}} \quad i = 1, \dots, n.$$

Standardized Residuals

The standardized residuals are defined by

$$r_i = \frac{\hat{\varepsilon}_i}{\hat{\sigma}\sqrt{1-h_{ii}}},$$

where h_{ii} is the *i*th diagonal element of the hat matrix.

Studentized Residuals

The studentized residuals are defined by

$$r_i^* = \frac{\hat{\varepsilon}_{(i)}}{\hat{\sigma}_{(i)}(1 + \mathbf{x}_i'(\mathbf{X}_{(i)}'\mathbf{X}_{(i)})^{-1}\mathbf{x}_i)^{1/2}} = \frac{\hat{\varepsilon}_i}{\hat{\sigma}_{(i)}\sqrt{1 - h_{ii}}} = r_i \left(\frac{n - p - 1}{n - p - r_i^2}\right)^{1/2}.$$

The studentized residuals are used to verify model assumptions and to discover outliers (see Sect. 3.4.4).

Box from our text book: Fahrmeir et al (2013): Regression. Springer. (p.126)

Simulating data and checking residuals

```
n=1000
beta=matrix(c(0,1,1/2,1/3),ncol=1)
set.seed(123)
x1=rnorm(n,0,1); x2=rnorm(n,0,2); x3=rnorm(n,0,3)
X=cbind(rep(1,n),x1,x2,x3)
y=X%*\%beta+rnorm(n,0,2)
fit=lm(y^x1+x2+x3)
yhat=predict(fit)
summary(fit)
ehat=residuals(fit); estand=rstandard(fit); estud=rstudent(fit)
plot(yhat,ehat,pch=20)
points(yhat,estand,pch=20,col=2)
#points(yhat,estud,pch=20,col=5)
```



Black: raw residuals, red: standardized residuals (identical to studentized here).

Examination of model assumptions

- 1. Linearity of covariates: $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
- 2. Homoscedastic error variance: $Cov(\varepsilon) = \sigma^2 I$.
- **3**. Uncorrelated errors: $Cov(\varepsilon_i, \varepsilon_j) = 0$.
- 4. Additivity of errors: $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
- 5. Assumption of normality: $\varepsilon \sim N_n(0, \sigma^2 I)$

Plotting residuals

- 1. Plot the residuals, r_i^* against the predicted values, \hat{y}_i .
 - Dependence of the residuals on the predicted value: wrong regression model?
 - Nonconstant variance: transformation or weighted least squares is needed?
- 2. Plot the residuals, r_i^* , against predictor variable or functions of predictor variables. Trend suggest that transformation of the predictors or more terms are needed in the regression.
- 3. Assessing normality of errors: QQ-plots and histograms of residuals. As an additional aid a test for normality can be used, but must be interpreted with caution since for small sample sizes the test is not very powerful and for large sample sizes even very small deviances from normality will be labelled as significant.
- 4. Plot the residuals, r_i^* , versus time or collection order (if possible). Look for dependence or autocorrelation.

Volume of a tree

Data for 31 trees of a certain kind in a national park in the US are given below. Three variables are measured for each tree. These are:

- D: The diameter of the tree measured in inches 1.5 m above ground level
- \blacktriangleright *H*: The height of the tree measured in feet.
- ► *V*: The volume of the tree measured in cubic feet.

Obs.	D	Н	V	Obs.	D	Н	V
1	8.3	70	10.3	17	12.9	85	33.8
2	8.6	65	10.3	18	13.3	86	27.4
3	8.8	63	10.2	19	13.7	71	25.7
4	10.5	72	16.4	20	13.8	64	24.9
5	10.7	81	18.8	21	14.0	78	34.5
6	10.8	83	19.7	22	14.2	80	31.7
7	11.0	66	15.6	23	14.5	74	36.3
8	11.0	75	18.2	24	16.0	72	38.3
9	11.1	80	22.6	25	16.3	77	42.6
10	11.2	75	19.9	26	17.3	81	55.4
11	11.3	79	24.2	27	17.5	82	55.7
12	11.4	76	21.0	28	17.9	80	58.3
13	11.4	76	21.4	29	18.0	80	51.5
14	11.7	69	21.3	30	18.0	80	51.0
15	12.0	75	19.1	31	20.6	87	77.0
16	12.9	74	22.2				

Volume of a tree

- If one wants to measure the volume of a tree the tree has to be cut down.
- But, height and diameter can be measured without cutting down the tree.
- Of interest: develop a model that can be used to estimate the tree volume from the height and diameter.

As an illustration assume we want to fit a linear model with V as response and D and H as covariates. What is the R^2 of this model?

Comment: if we start with the volume of a cylinder (area of circle times height) we may suggest a different regression model (on the log scale). Which model?

Volume: height and diameter

```
fit <- lm(Volume~.,data=ds)
summary(fit)</pre>
```

Coefficients:

Estimate Std. Error t value Pr(>|t|) (Intercept) -57.9877 8.6382 -6.713 2.75e-07 *** Diameter 4.7082 0.2643 17.816 < 2e-16 *** Height 0.3393 0.1302 2.607 0.0145 * ---Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '

Residual standard error: 3.882 on 28 degrees of freedom Multiple R-squared: 0.948,Adjusted R-squared: 0.9442 F-statistic: 255 on 2 and 28 DF, p-value: < 2.2e-16

Volume of a tree: IQ of lumberjack added

- We want to add the IQ of the lumberjack that cut down the tree as a covariate in the model.
- This should for obvious reasons not be a good predictor for the volume of the tree.
- To mimic this situation we simulate new data to resemble the IQ of different lumberjacks by drawing data from the normal distribution with mean 100 and standard deviation 16, and since we have 31 trees we simulate 31 observations.
- Q: will the R² of this new model be higher than the R² of the previous model?

Volume: height and diameter – and IQ of lumberjack

set.seed(123) # reproducible results
iq <- rnorm(31,100,16)
fit2 <- lm(Volume~Height+Diameter+iq,data=ds)
summary(fit2)</pre>

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	-61.03399	10.20868	-5.979	2.24e-06	***
Height	0.34099	0.13176	2.588	0.0154	*
Diameter	4.72507	0.26906	17.561	2.68e-16	***
iq	0.02704	0.04678	0.578	0.5681	
Signif. cod	es: 0 '***	' 0.001 '**	· 0.01	** 0.05	'.' 0.1

Residual standard error: 3.929 on 27 degrees of freedom Multiple R-squared: 0.9486,Adjusted R-squared: 0.9429 F-statistic: 166.1 on 3 and 27 DF, p-value: < 2.2e-16

Acid rain in Norwegian lakes

Data on n = 26 lakes, with

- ► y: measured pH in lake,
- ► x1: *SO*₄: sulfate (the salt of sulfuric acid),
- ► x2: N0₃: nitrate (the conjugate base of nitric acid),
- ► x3: *Ca*: calsium,
- ► x4: latent *AI*: aluminium,
- x5: organic substance,
- ► x6: area of lake,
- ► x7: position of lake (Telemark or Trøndelag),

We would like to use a regression model with pH of the lake as the response. Should we fit a model will all 7 covariates, or choose a subset?

Simulated data (Fahrmeir et al: Fig 3.17)

True model:

$$Y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \varepsilon_i$$

Known that the model is polynomial in nature, but not up to which degree.

Try to fit polynomial also with higher order terms.

New: in addition to the data set to be used to fit the regression (called *training set*) also a data set to assess the model fit is present (called a *validation* set).

Mean Squared Error (MSE) is a scaled version of the SSE, that is $\frac{1}{n} \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2$.



Fig. 3.17 Simulated training data y_i [panel (**a**)] and validation data y_i^* [panel (**b**)] based on 50 design points x_i , i = 1, ..., 50. The true model used for simulation is $y_i = -1 + 0.3x_i + 0.4x_i^2 - 0.8x_i^3 + \varepsilon_i$ with $\varepsilon_i \sim N(0, 0.07^2)$. Panels (**c**-**e**) show estimated polynomials of degree l = 1, 2, 5 based on the training set. Panel (**f**) displays the mean squared error MSE(*l*) of the fitted values in relation to the polynomial degree (*solid line*). The *dashed line* shows MSE(*l*), if the estimated polynomials are used to predict the validation data y_i^*

Figure from our text book: Fahrmeir et al (2013): Regression. Springer. (p.140)

Simulated data (Fahrmeir et al: Fig 3.18, Tab3.3, Tab3.4)

True model:

$$Y \sim N(-1+0.3x_1+0.2x_3, 0.2^2)$$

where also $x_2 = x_1 + u$ is observed ($u \sim$ uniform in 0,1). The variables x_1 and x_3 are uncorrelated.



ig. 3.18 Scatter plot matrix for the variables y, x_1 , x_2 , and x_3

Figure from our text book: Fahrmeir et al (2013): Regression. Springer. (p.141)

Variable	Coefficient	Standard error	t-value	p-value	95 % Con	fidence interval
intercept	-0.970	0.047	-20.46	< 0.001	-1.064	-0.877
x_1	0.146	0.187	0.78	0.436	-0.224	0.516
<i>x</i> ₂	0.027	0.177	0.15	0.880	-0.323	0.377
<i>x</i> ₃	0.227	0.052	4.32	< 0.001	0.123	0.331

Table 3.3 Results for the model based on covariates x_1 , x_2 , and x_3

Table 3.4 Results for the correctly specified model based on covariates x_1 and x_3

Variable	Coefficient	Standard error	t-value	p-value	95 % Con	fidence interval
intercept	-0.967	0.039	-24.91	< 0.001	-1.042	-0.889
x_1	0.173	0.055	3.17	0.002	0.065	0.281
<i>x</i> ₃	0.226	0.052	4.33	< 0.001	0.123	0.330

Table from our text book: Fahrmeir et al (2013): Regression. Springer. (p.142)

Irrelevant and/or missing covariates in the regression

Irrelevant : variables that are included in the regression but should not have been.

missing : variables that are not included, but should have been.

Classical linear model with identically normally distributed random errors, $Cov(\varepsilon) = \sigma^2 I$, but now look at misspecification of $E(\mathbf{Y})$. Suppose that the true model is

$$\begin{aligned} \mathbf{Y} &= \mathbf{X}_1 \boldsymbol{\beta}_1 + \mathbf{X}_2 \boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}, \\ \boldsymbol{\varepsilon} &\sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}), \end{aligned} \tag{1}$$

where we have partitioned the design matrix into two parts X_1 $(n \times p_1)$ and X_2 $(n \times p_2)$ and β_1 and β_2 are unknown p_1 - and p_2 -dimensional vectors of regression coefficients $(p = p_1 + p_2)$. Assume that we ignore the covariates in X_2 and fit the model

$$oldsymbol{Y} = oldsymbol{X}_1 lpha_1 + \delta, \ \delta \sim N_n(\mathbf{0}, \tau^2 oldsymbol{I}).$$
 (2)

Here α_1 is used in place of β_1 to emphasize that α_1 (and estimates thereof) will in general be different from β_1 in the true model. The least squares estimator for model (2) is $\hat{\alpha_1} = (\boldsymbol{X}_1^T \boldsymbol{X}_1)^{-1} \boldsymbol{X}_1^T \boldsymbol{Y}.$

Two subsets of covariates (cont.)

Find the expected value and covariance matrix of $\hat{\alpha_1}$ under the true model.

$$E(\hat{\boldsymbol{\alpha}_1}) = \boldsymbol{\beta}_1 + (\boldsymbol{X}_1^T \boldsymbol{X}_1)^{-1} \boldsymbol{X}_1^T \boldsymbol{X}_2 \boldsymbol{\beta}_2$$

We see that the bias term for $\hat{\alpha}_1$ is $(\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{X}_2 \beta_2$. When is the bias term equal to zero?

$$\operatorname{Cov}(\hat{\boldsymbol{\alpha}_1}) = \sigma^2 (\boldsymbol{X}_1^T \boldsymbol{X}_1)^{-1}$$

Observe, $\operatorname{Cov}(\hat{\alpha_1})$ is not dependent on β_2 .

Missing covariates: findings

- Bias : The estimator for the (true) covariates (in the model) is only unbiased if the true and missing covariates are uncorrelated (orthogonal design) in the data.
- Variance : The variance of the estimator for the true covariates may be smaller based on the model with the missing covariates (than for the correctly specified model), and even the sum of the bias² and the variance may better for the model with the missing variables. So the sparse model may be better on overall (even though it is biased).

Irrelevant covariates included: findings

- Bias : The estimator for the true covariates are unbiased, also if irrelevant covariates are included.
- Variance : The model with the irrelevant covariants have larger variance for the true covariates, compared with the model without the irrelevant covariates. So, again sparse model is the best.

Irrelevant and/or missing covariates in the regression

- Irrelevant : variables that are included in the regression but should not have been.
 - missing : variables that are not included, but should have been.

Conclusion in book: the model should not contain irrelevant covariates, and we should aim for a sparse model.

Law of parsimony

If two models are not very different – then always choose the simplest one
Today

- T-test for significance of one regression coefficient.
- Residuals: standardized (or studentized) preferred.
- Significance of regression based on F-test with SSR/(p-1) divided by SST/(n-1).
- R² gives the proportion of variability explained by the regression model.

$$R^2 = rac{\text{SSR}}{\text{SST}} = 1 - rac{\text{SSE}}{\text{SST}}$$

and will never decrease if new covariates are added to the model.

Model selection: want to choose the model that minimize the expected squared prediction error.

Previously:
$$Y = X\beta + E$$
, $E \sim Nn(0, \sigma^2 I)$
 $\hat{\beta}_j$ is the jth element of $\hat{\beta} = (X^T X)^{-1} X^T Y$
 $\hat{\beta} \sim Np(\beta, (X^T X)^{-1} \sigma^2)$
 $\hat{\sigma}^2 = n-p \quad \hat{\epsilon}^T \hat{\epsilon} = \frac{SSE}{n-p}$
 $\hat{\epsilon} = Y - \hat{Y} = Y - X\hat{\rho}$



These for association (linear) between response Y and x_j : Ho: $B_j = 0$ is $H_i: B_j \neq 0$ When Ho is true: $T_j o = \frac{B_j - 0}{V_{cjj}!} \sim t_{n-p}$ Here statistic

p-value:
$$P(||T_{i0}| \ge |t_{i0}||_{1} \text{ Ho hrw})$$

= $2 \cdot P(|T_{i0}| \ge |t_{i0}||_{1} \text{ Ho hrw})$
Reject Ho when $|t_{i0}| \ge t_{x_{2}} \text{ prop}$ Sign. level d.
Ex: Acid rain : lineer association between
Soy and pH
Ho: $B_{1} = 0$ vs H_{1} : $B_{1} \neq 0$
From summery of lm in R (slide)
 $t_{10} = -5,362$
 $p_{1} = 15$
 $p_{1} = 15$
 $p_{1} = 15$
 $p_{1} = 15$
 $p_{1} = 15,362$
 $p_{1} = 15,362$
 $p_{2} = \frac{4.5 \cdot 10^{5}}{1}$
 $R: 2*(A pt(5.3b2, 18))$
Reject Ho fir all area to left

Residuals (agein)

$$\hat{\mathcal{E}} = Y - \hat{Y}, \quad \hat{\mathcal{E}} \sim N_n(0, 0^{\circ}(\mathbf{I} - \mathbf{H}))$$

R: residuals (kt)

The residuate have heteroscedestic verience

 $Var(\hat{\mathcal{E}}) = 0^{\circ}(1 - hii)$ and $Gv(\hat{\mathcal{E}}_i, \hat{\mathcal{E}}_j) = \sigma^{\circ}(0 - hij)$

Cen in general bea

 $\neq 0, but in nost$

some experence shows that

 $\hat{\mathcal{E}} = \frac{\hat{\mathcal{E}}_i}{\hat{\mathcal{E}} \vee 1 - hii}$ will be (approx.)

 $r_2 = \frac{\hat{\mathcal{E}}_i}{\hat{\mathcal{E}} \vee 1 - hii}$ homoscedestic.

R: rstanderd (Tit)

Stadentized residuals:

 $r_i^* = see slide$

Hee studentized!

See expemple on slide

 $R: rstudent(pt)$

 $for r_i vs \hat{\mathcal{E}}_i$

Analysis of Verience decomposition and
$$R^2$$

 y_1, \dots, y_n and $\tilde{y} = \pi \sum_{i=1}^{n} y_i$
 $\sum_{i=1}^{n} (y_i - \tilde{y})^2 = \sum_{i=1}^{n} (y_i - \hat{y}_i + \hat{y}_i - \tilde{y})^2 = \dots =$
Sums of squared $= \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \sum_{i=1}^{n} (\hat{y}_i - \tilde{y})^2$
SST
Sums of squared sums of squared sums of sq
error regression
 BSE
SSR
 y explained
with vectors and matrices.
 $\xi \tau \in$
 $(T - \pi MT)Y = Y^T (I - H)Y + Y^T (H - \pi MT)Y$
SST
 SSE
 SSE
 SSR

This is used to define:

Is the regression significent?
Ho:
$$B_1 = B_2 = \dots = \beta_k = 0$$
 US
 $M_1: at least one B_j \neq 0$ $j=1,\dots,k$
(SST-SSE)
Test statistico:
 $F = \frac{SSR/k}{SSE/(n-p)} \sim F_k, n-p$
 $Prove this in Part 3 in
a general setting$

Ex: Acidrein:

F-observed: 34.15
$$P(F_{7,18} > 34.15) = 3.9.0^{-7}$$

 $k = 7, n-p = 18$ $P(F_{7,18} > 34.15) = 3.9.0^{-7}$
 $25 8$ p -value

Big model:
$$Y_i = potpix_i + p_2 x_{2i} + p_{3}x_{1} + e_i$$

Small model: $Y_i = potp_1 x_{1i} + p_2 x_{2i} + e_i$
 $R^{e}_{oig} \ge R^{2}_{onall} \iff since p_{3}$ is found
to minimize $SSE = \hat{e}^{\dagger}\hat{e}$, thus maximize
 $R^{e} = 1 - \frac{SSE}{SST}$
 $R^{2}B_{ig} = R^{2}snall ef \hat{p}_{3} = 0$
if $p_{3}^{2} \neq 0$ the $SSE_{Big} < SSE_{snall} = rod$
 $R^{e}_{big} \ge R^{2}_{onall}$.

R^e will always increase (or stay unchanged) when a new coversite is added to the rodel.

Next: more on choosing a good model, and
then
$$\frac{R^{2}}{adj} = 1 - \frac{SSE/(n-p)}{SST/(n-1)} e^{penolizing} adding$$
is one criterion to use instead of R²
for model selection.

Model choice end verieble selection (F3.Y]

Question 1: Is a full model (all available
avanation filted) the best model?
good interpretability good for future predictions
Deta is divided into
Training Validetian
set
T available
parameter model walkable
parameter model fit
MSE =
$$n \sum_{i=1}^{n} (Y_i - Y_i)^2$$

Now calculable the MSE on the training end on the
validation of
model an plat as afunction of
model an plat is a afunction of
model an plat is a afunction of
model an plat is a afunction of
MSE = $n \sum_{i=1}^{n} (Y_i - Y_i)^2$
Now calculable the MSE on the training end on the
validation of
model an plat is a afunction of
model an platity:
Malideton
Malideton
Malideton
Malideton in model

Answer 1: No, this may lead to overlitting = fitting the trend + the norse !

=> so, whet cen we do instead?

TMA4267 Linear Statistical Models V2017 (L12)

Part 2: Linear regression: Model selection [F:3.4] Transformation and Taylor expansion Quiz

Mette Langaas

Department of Mathematical Sciences, NTNU

To be lectured: February 24, 2017

What is the "best" model?

Acid rain in Norwegian lakes, data on n = 26 lakes, with

- ► y: measured pH in lake,
- ► x1: SO₄: sulfate (the salt of sulfuric acid),
- ► x2: N0₃: nitrate (the conjugate base of nitric acid),
- ► x3: *Ca*: calsium,
- ► x4: latent *AI*: aluminium,
- ► x5: organic substance,
- ► x6: area of lake,
- x7: position of lake (Telemark or Trøndelag),

Topic: choosing the "best" linear regression model!

- First, debunk popular strategies (based on simulations studies were we knew the "true" model):
 - Popular 1: fit all available covariates.
 Problem: overfitting (=fitting trends and noise).
 - Popular 2: fit all available covariates, then remove the insignificant ones (=those β_j where H₀ : β_j = 0 is not rejected).

Simulated data (Fahrmeir et al: Fig 3.18, Tab3.3, Tab3.4)

True model:

$$Y \sim N(-1 + 0.3x_1 + 0.2x_3, 0.2^2)$$

where also $x_2 = x_1 + u$ is observed ($u \sim$ uniform in 0,1). The variables x_1 and x_3 are uncorrelated.



ig. 3.18 Scatter plot matrix for the variables y, x_1 , x_2 , and x_3

Figure from our text book: Fahrmeir et al (2013): Regression. Springer. (p.141)

Variable	Coefficient	Standard error	t-value	p-value	95 % Confidence interval		
intercept	-0.970	0.047	-20.46	< 0.001	-1.064	-0.877	
x_1	0.146	0.187	0.78	0.436	-0.224	0.516	
<i>x</i> ₂	0.027	0.177	0.15	0.880	-0.323	0.377	
<i>x</i> ₃	0.227	0.052	4.32	< 0.001	0.123	0.331	

Table 3.3 Results for the model based on covariates x_1 , x_2 , and x_3

Table 3.4 Results for the correctly specified model based on covariates x_1 and x_3

Variable	Coefficient	Standard error	t-value	p-value	95 % Confidence interval		
intercept	-0.967	0.039	-24.91	< 0.001	-1.042	-0.889	
x_1	0.173	0.055	3.17	0.002	0.065	0.281	
<i>x</i> ₃	0.226	0.052	4.33	< 0.001	0.123	0.330	

Table from our text book: Fahrmeir et al (2013): Regression. Springer. (p.142)

Topic: choosing the "best" linear regression model!

- First, debunk popular strategies (based on simulations studies were we knew the "true" model):
 - Popular 1: fit all available covariates.
 Problem: overfitting (=fitting trends and noise).
 - Popular 2: fit all available covariates, then remove the insignificant ones (=those β_j where H₀ : β_j = 0 is rejected). Problem: may also remove important covariates that are correlated with unimportant ones but insignificant because being masked by the unimportant ones.

Study of irrelevant and missing covariates:

- Irrelevant : variables that are included in the regression but should not have been (IQ of lumberjack)
- missing : variables that are not included, but should have been (omitting height in the tree volum example)

Conclusion in book: the model should not contain irrelevant covariates, and we should aim for a sparse model.

Take home message is the "Law of parsimony": *If two models* are not very different – then always choose the simplest one.

A model is a simplification or approximation of reality and hence will not reflect all of reality.

George Box noted that "all models are wrong, but some are useful". While a model can never be "truth"a model might be ranked from very useful, to useful, to somewhat useful to, finally, essentially useless.

Burnham, K. P.; Anderson, D. R. (2002), Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach.

Expected squared prediction error (SPSE)

Possible criterion we want to minimize: SPSE. Definition (j, M, ... given in classnotes)

$$\mathsf{SPSE} = \sum_{j=1}^{J} \mathrm{E}((Y_j - \hat{Y}_{jM})^2)$$

can be written as:

$$SPSE = \sum_{j=1}^{J} E((Y_j - \hat{Y}_{jM})^2) = n\sigma^2 + |M|\sigma^2 + \sum_{j=1}^{J} (\mu_{jM} - \mu_j)^2$$

Problem: Not useful on practise since μ_j and σ^2 are unknown. Plan: Find a way to estimate SPSE and then choose the model M with the minimum SPSE!

How to estimate SPSE?

$$\mathsf{SPSE} = \sum_{j=1}^J \mathrm{E}((Y_j - \hat{Y}_{jM})^2)$$

Assume we have fitted a model M with |M| regression parameters.

1. Use new (independent) data – if available (seldom the case):

$$\widehat{SPSE} = \sum_{j=1}^{J} (Y_j - \hat{Y}_{jM})^2$$

Cross-validation: mimic new data by dividing data into k folds (popular is k = n and k = 10). In a for-loop let j = 1,..., k, and use all folds except fold j to estimate regression parameter, and use the jth fold to calculated the SPSE. Sum across folds.
 Choose the model M that minimizes the SPSE.

Cross-validation (5-fold)



Will be taught in TMA4300 Computational statistics and will be a backbone in TMA4268 Statistical Learning. http://blog-test.goldenhelix.com/wp-content/uploads/2015/04/B-fig-1.jpg

How to estimate SPSE?

$$\mathsf{SPSE} = \sum_{j=1}^J \mathrm{E}((Y_j - \hat{Y}_{jM})^2)$$

Assume we have fitted a model M with |M| regression parameters.

3. Use existing data (only): It can be shown that

 $E(\widehat{SPSE}) = SPSE - 2 | M | \sigma^2$ when used on the same data that was used to make the prediction, so a better estimate for existing data is

$$\widehat{SPSE} = \sum_{i=1}^{n} (Y_i - \hat{Y}_{iM})^2 + 2|M|\hat{\sigma}^2 = SSE + 2|M|\hat{\sigma}^2$$

where $\hat{\sigma}^2$ is the same for all models M, and is often estimated using the most complex model under study.

4. Other criteria: all have the same form; a first term based on SSE (or R^2) for model M, and a second term penalizing the model complexity.

Choose the model M that minimizes the \widehat{SPSE} .

For models with the same model complexity – easy solution: SSE

- Estimators for SPSE to be used on the same data as to be used for estimating the model parameters have the same form; a first term based on SSE (or R^2) for model M, and a second term penalizing the model complexity.
- If we consider two models with the same model complexity then SSE can be used to choose between these models.

Acid rain (1). Best subset

For 1,...,7 covariates: fit all possible models, and report the model with the smallest SSE (given below) for each value for the model complexity. Explain what you see! How many models have been searched for each model complexity?

```
regfit.full=regsubsets(y~.,data=ds)
sumreg <- summary(regfit.full)
Subset selection object
Call: regsubsets.formula(y ~ ., data = ds)
Selection Algorithm: exhaustive</pre>
```

Names: x1: SO_4 , x2: NO_3 , x3: Ca, x4: latent AI, x5: organic substance, x6: area of lake, x7: position of lake (Telemark or Trøndelag).

 R^2 adjusted (corrected) Mallows' C_p Akaike Information Criterion (AIC) Bayesian Information Criterion (BIC)

NB: there is no overall best choice for criterion - all of these are used.

 R^2 adjusted (corrected)

 \hat{Y}_i is from fitting the regression model M. Remember, for a regression model (with intercept) we have the SST=SSR+SSE.

$$SST = \sum_{i=1}^{n} (Y_i - \bar{Y})^2$$

$$SSE = \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2$$

$$R^2 = 1 - \frac{SSE}{SST}$$

$$R^2_{adj} = 1 - \frac{\frac{SSE}{n-p}}{\frac{SST}{n-1}} = 1 - \frac{n-1}{n-p}(1-R^2)$$

Choose the model with the *largest* R_{adj}^2 .

"All" statistical software outputs this automatically! However, Fahrmeir et al (2013) believes that the penalty n - p is too small.

Happiness (n = 39)

Are love and work the important factors determining happiness?

- y, happiness. 10-point scale, with 1 representing a suicidal state, 5 representing a feeling of «just muddling along», and 10 representing a euphoric state.
- > x_1 , money. Annual family income in thousands of dollars.
- x₂, sex. Sex was measured as the values 0 or 1, with 1 indicating a satisfactory level of sexual activity.
- x₃, love. 3-point scale, with 1 representing loneliness and isolation, 2 representing a set of secure relationships, and 3 representing a deep feeling of belonging and caring in the context of some family or community.
- x₄, work. 5-point scale, with 1 indicating that an individual is seeking other employment, 3 indicating the job is OK, and 5 indicating that the job is enjoyable.

Data taken from library faraway, data set happy.

Нарру

```
> allreg=regsubsets(happy~.,data=happy)
> sumreg <- summary(allreg)</pre>
> sumreg
Subset selection object
Call: regsubsets.formula(happy ~ ., data = happy)
1 subsets of each size up to 4
Selection Algorithm: exhaustive
        money sex love work
1 (1) " " " " " " "
2 (1) " " " " " * " * "
3 (1) "*" " "*" "*"
4 (1) "*" "*" "*"
```

	money	sex	love	work	N	р	R^2	$R_{\rm adj}^2$
1	0.014				1	0.000747	7.3	4.8
2		-0.130			1	1	0.1	-2.6
3			2.270		1	8.35e-24	61.5	60.5
4				0.990	1	1.36e-13	29.1	27.2
5	0.016	-0.508			2	0.0504	8.8	3.8
6	0.009		2.206		2	8.77e-19	64.5	62.5
7	0.012			0.961	2	3.68e-10	34.6	31.0
8		-0.277	2.279		2	5.55e-18	62.0	59.9
9		0.610		1.079	2	3.48e-09	31.2	27.4
10			1.959	0.511	2	5.75e-20	68.1	66.3
11	0.011	-0.536	2.209		3	9.49e-16	66.2	63.3
12	0.011	0.305		1.009	3	1.84e-07	35.1	29.5
13	0.009		1.902	0.504	3	2.63e-17	70.9	68.4
14		0.108	1.944	0.530	3	2.22e-16	68.1	65.4
15	0.010	-0.149	1.919	0.476	4	9.89e-15	71.0	67.6

Intercept included, N = p - 1, *p*-value for significance of regression. $R^2 = 1 - \frac{SSE}{SST}$, $R_{adj}^2 = 1 - \frac{\frac{SSE}{n-p}}{\frac{SST}{n-1}}$. Which model to prefer?



 \hat{Y}_i is from fitting regression model M. Mallows is the name of a person.

$$C_{p} = \frac{\sum_{i=1}^{n} (Y_{i} - \hat{Y}_{i})^{2}}{\hat{\sigma}^{2}} - n + 2|M|$$

Minimizing Cp gives the same optimal model as minimizing SPSE.

See Exam V2015 Problem 3 for an in depth explanation of the theory behind Mallow's *Cp*.

Akaike information criterion – one of the most widely used. Designed for likelihood-based inference.

For a normal regression model:

$$AIC = n\ln(\hat{\sigma}^2) + 2(|M| + 1)$$

Choose the model with the minimum AIC.

Bayesian information criterion.

For a normal regression model:

$$\mathsf{BIC} = n \ln(\hat{\sigma}^2) + \ln(n)(|M| + 1)$$

Choose the model with the minimum BIC.

AIC and BIC are motivated in very different ways, but the final result for the normal regression model is very similar.

BIC has a larger penalty than AIC $(\log(n)vs.2)$, and will often give a smaller model (=more parsimonious models) than AIC.

Happy: Mallows' C_p



Happy: BIC



Acid rain (2)

Call: regsubsets.formula(y ~ ., data = ds)
1 subsets of each size up to 7
Selection Algorithm: exhaustive

x1 x2 x3 x4 x5 x6 x7 1 (1) " " " " " " " " " " " " " " (1) "*" " " "*" " " " " " " " " " 2 3 (1) "*" "*" "*" " " " " " " " " 4 (1) "*" "*" "*" " " " " " " 5 (1) "*" "*" " " "*" " "*" (1) "*" "*" "*" "*" "*" " "*" 6 (1) "*" "*" "*" "*" "*" "*" "*" 7 # to mimic test set: which.max(sumreg\$adjr2) #5 which.min(sumreg\$cp) #3 which.min(sumreg\$bic) #3 # so, model 3 or 5 is suggested for us # model 3: x1+x2+x3 # model 5: x1+x2+x3+x5+x7

Acid rain, BIC,



Practical use of the model criteria

- All subset selection: use smart "leaps and bounds" algorithm, works fine for number of covariates in the order of 40.
- Forward selection: choose starting model (only intercept), then add one new variable at each step - selected to make the best improvement in the model selection criteria. End when no improvement is made.
- Backward elimination: : choose starting model (full model), then remove one new variable at each step - selected to make the best improvement in the model selection criteria. End when no improvement is made.
- Stepwise selection: combine forward and backward.
Acid rain (3): stepAIC

```
> all=lm(happy~.,data=happy)
> stepAIC(all)
Start: AIC=9.08
happy ~ money + sex + love + work
       Df Sum of Sq
                      RSS
                             AIC
              0.142 38.229 7.221
       1
- sex
<none>
                    38.087 9.076
- money 1
           3.782 41.869 10.768
- work
      1
           6.386 44.473 13.122
- love 1 47.272 85.359 38.549
Step: AIC=7.22
happy ~ money + love + work
       Df Sum of Sq
                      RSS
                             AIC
                    38.229 7.221
<none>
           3.723 41.952 8.846
- money 1
- work 1 8.410 46.639 12.976
- love 1 47.742 85.971 36.828
Call:
lm(formula = happy ~ money + love + work, data = happy)
Coefficients:
(Intercept)
                               love
                 money
 -0.185936
               0.008959
                           1.901709
                                       0.503602
```

work

Acid rain (4): Forward

regfitF=regsubsets(y~.,data=ds,method="forward")
sumregF <- summary(regfitF)
Selection Algorithm: forward</pre>

				x1	x2	xЗ	x4	x5	x6	x7
1	(1)	11 11	11 11	11 11	"*"	11 11	11 11	11 11
2	(1)	11 11	11 11	"*"	"*"	11 11	11 11	11 11
3	(1)	"*"	"*"	"*"	11 11	11 11	11 11	11 11
4	(1)	"*"	"*"	"*"	"*"	11 11	11 11	11 11
5	(1)	"*"	"*"	"*"	"*"	"*"	11 11	11 11
6	(1)	"*"	"*"	"*"	"*"	"*"	11 11	"*"
7	(1)	"*"	"*"	"*"	"*"	"*"	"*"	"*"
which.max(sumregF\$adjr2)#5										
which.min(sumregF\$cp) #3										
which.min(sumregF\$bic) #3										

Acid rain (5): Backward

regfitB=regsubsets(y~.,data=ds,method="backward")
sumregB <- summary(regfitB)
Selection Algorithm: backward</pre>

				x1	x2	xЗ	x4	x5	x6	x7
1	(1)	11 11	11 11	"*"	11 11	11 11	11 11	11 11
2	(1)	"*"	11 11	"*"	11 11	11 11	11 11	11 11
3	(1)	"*"	"*"	"*"	11 11	11 11	11 11	11 11
4	(1)	"*"	"*"	"*"	11 11	"*"	11 11	11 11
5	(1)	"*"	"*"	"*"	11 11	"*"	11 11	"*"
6	(1)	"*"	"*"	"*"	"*"	"*"	11 11	"*"
7	(1)	"*"	"*"	"*"	"*"	"*"	"*"	"*"
which.max(sumregB\$adjr)#5										
<pre># backward finds same as best subset</pre>										
which.min(sumregB\$cp) #3										

- Influential observations and outliers: impact of specific observations on model fit.
- Collinearity analysis: Highly correlated variables cause imprecise estimation of the regression parameters. (Why? Look at diagonal elements of Cov(β) = σ²(X^TX)⁻¹, and look back to Problem 2 in the start of this lecture.)
- Examination of model assumptions: residual plots!

Influential observations- and outliers

- Observations that significantly affect inferences drawn from the data are said to be influential.
- ► The leverage, h_{ii}, associated with the *i*th datapoint measures "how far the *i*th observation is from the other n − 1 observations".
- Methods for assessing influential observations may be be based on change in *β* estimate when observations are deleted.
- Always investigate possible causes of an influential observation (if possible).
- Cook's distance can be used to identify influential observations.
- Robust methods (median, quantile regression) can be useful.

Want to understand more? Read for yourself in Fahrmeir et al (2013): p 160-166.

Transformations

- Multiplicative or additive model?
- Box–Cox transform with profile likelihood.
- Stabilizing the variance.

Galapagos islands, Model A, Exam V2014 Problem 2



Normal Q–Q Plot

Box–Cox plot



Box–Cox transformation plot based on Model A for the Galapagos data set, RecEx4. Line at x = 1/3.

Galapagos islands, Model B, Exam V2014 Problem 2



Normal Q–Q Plot

Approximation of E and Var for nonlinear functions

- Have RV X, with mean $E(X) = \mu$ and some variance Var(X).
- Want to look at a nonlinear function of X, called g(X).
- Aim: find an approximation to E(g(X)) and Var(g(X)).
- And, the same for two RVs X_1 and X_2 with $g(X_1, X_2)$.

Looking at residual plots from a regression model the conclusion was to analyse data of *BMI* on the natural logarithmic scale. After a regression model was fitted the predicted value for the ln(BMI) for a specific combination of the covariates was found to be 3.2151 with an estimated standard deviation of 0.1656. Use approximate methods to arrive at an estimate of the predicted value and estimated standard deviation on the original scale, kg/m², and not on the logarithmic scale.

E(g(X) and Var(g(X)))

- Let g(X) be a general function. When is E(g(X)) = g(E(X))?
 - When g(X) is a linear function of X.
- What can we do if this is not the case?
 - We can calculate $E(g(X)) = \int_{-\infty}^{\infty} g(x)f(x)dx$ when X is continuous, or a version thereof in the discrete case,
 - or if g is monotone we can use the transformations formula to find the distribution of Y = g(X) and then calculate E(Y) and Var(Y), if possible.
- What if we only know E(X) = μ and Var(X) = σ² and not f(x)?
 - ► Use a Taylor series approximation of g(X) around g(µ). g need to be differentiable.

First order Taylor approximation of g(X) around μ .

$$g(X) \approx g(\mu) + g'(\mu)(X - \mu)$$

This leads to the following approximations:

$$\mathrm{E}(g(X)) \approx g(\mu)$$

 $\mathrm{Var}(g(X)) \approx [g'(\mu)]^2 \mathrm{Var}(X)$

Treatment of tennis elbow (exam TMA4255 V2012, 3b)

The term *tennis elbow* is used to describe a state of inflammation in the elbow, causing pain. This injury is common in people who play racquet sports, however, any activity that involves repetitive twisting of the wrist (like using a screwdriver) can lead to this condition. The condition may also be due to constant computer keyboard and mouse use.

In a randomized clinical study the aim was to compare three different methods for treatment of tennis elbow,

- A: physiotherapy intervention,
- B: corticosteroid injections and
- C: wait-and-see (the patients in the wait-and-see group did not get any treatment but was told to use the elbow as little as possible).

We will look at the short-term effect of treatment by studying measurements at 6 weeks. All patients participating in the study only had one affected arm.

We will look at the outcome measure called *pain-free grip force*. This was measured by a digital grip dynamometer and normalized to the grip force of the unaffected arm. A pain-free grip force of 100 would mean that the affected and the unaffected arm performed equally good. Summary statistics for each of the treatment groups.

Treatment	Sample size	Average	Standard deviation
A (physiotherapy)	63	70.2	25.4
B (injection)	65	83.6	22.9
C (wait-and-see)	60	51.8	23.0
Total	188	69.0	

Example 2: Exam TMA4255 V2012 3d (fraction)

Let μ_A be the expected pain-free grip force for a population where the physiotherapy intervention treatment is used to treat tennis elbow, and μ_C be the expected pain-free grip force for a population where the wait-and-see treatment is used. Define the relative difference between these two expected values as

$$\gamma = \frac{\mu_A - \mu_C}{\mu_C}.$$

This can be interpreted as the expected relative gain by using physiotherapy instead of wait-and-see. Based on two independent random samples of size n_A and n_C from the physiotherapy and wait-and-see treatment groups, respectively, suggest an estimator, $\hat{\gamma}$, for γ .

Use approximate methods to find the expected value and variance of this estimator, that is, $E(\hat{\gamma})$ and $Var(\hat{\gamma})$.

Bivariate function: first order Taylor

 X_1 is a RV with $\mu_1 = E(X_1)$ and X_2 is a RV with $\mu_2 = E(X_2)$. Let g be a bivariate function of X_1 and X_2 , and define

$$g_1'(\mu_1, \mu_2) = \frac{\partial g(x_1, x_2)}{\partial x_1} |_{x_1 = \mu_1, x_2 = \mu_2}$$
$$g_2'(\mu_1, \mu_2) = \frac{\partial g(x_1, x_2)}{\partial x_2} |_{x_1 = \mu_1, x_2 = \mu_2}$$

First order Taylor approximation:

 $g(X_1, X_2) \approx g(\mu_1, \mu_2) + g'_1(\mu_1, \mu_2)(X_1 - \mu_1) + g'_2(\mu_1, \mu_2)(X_2 - \mu_2)$

Bivariate function: first order Taylor

$$\begin{split} \mathrm{E}(g(X_1,X_2)) &\approx g(\mu_1,\mu_2) \\ \mathrm{Var}(g(X_1,X_2)) &\approx [g_1'(\mu_1,\mu_2)]^2 \mathrm{Var}(X_1) + [g_2'(\mu_1,\mu_2)]^2 \mathrm{Var}(X_2) + \\ & 2 \cdot g_1'(\mu_1,\mu_2) \cdot g_2'(\mu_1,\mu_2) \mathrm{Cov}(X_1,X_2) \end{split}$$

From Tabeller og formler i statistikk.

Rekkeutvikling

En første ordens Taylorutvikling av funksjonen $g(X_1, \ldots, X_n)$ omkring $g(\mu_1, \ldots, \mu_n)$, der $E(X_i) = \mu_i$, $i = 1, \ldots, n$, gir approksimasjonene

$$\mathbf{E}[g(X_1,\ldots,X_n)] \approx g(\mu_1,\ldots,\mu_n),$$

$$\mathbf{Var}[g(X_1,\ldots,X_n)] \approx \sum_{i=1}^n \left(\frac{\partial g(\mu_1,\ldots,\mu_n)}{\partial \mu_i}\right)^2 \mathbf{Var}(X_i) + 2\sum_{i>j} \frac{\partial g}{\partial \mu_i} \frac{\partial g}{\partial \mu_j} \mathbf{Cov}(X_i,X_j).$$

Today

- Choosing between models of equal model complexity: choose the model with the minimum SSE.
- Choosing between models of *different model complexity*: Model selection based on penalized criteria (Mallows Cp, R²_{adj},AIC and BIC). Try out on RecEx4 and Compulsory Exercise 2.
- BoxCox transformation: see RecEx4.
- Work for for yourself: Taylor solution to E and Var of nonlinear function, useful when you want to look at transformations of the data or functions of parameter estimates.

Summary of Part 2 in Kahoot!

F j=l E ((Y - Y H)²) = SPSE new obs predicted value besed on Bri 1



PLAN: elimple SPSE and choose the M that minimizes this estimate.

Finding the best model
$$\leftarrow$$
 all subsets method
indicat spece
1) Have k coverience that might be used
 $Y_i = \beta_0 + \beta_1 x_{1i} + \cdots + \beta_{1i} x_{ki} + \epsilon_i$
How many possible module can to make
(weat to have intracept)?
 $2 \cdot 2 \cdot 2 \cdots 2 = 2^{k}$ possible module
fit all possible 2^{k} module.
2) For $M = h 0, 1, 2, \cdots, k$ choose the
model with the smallert SSE.
Ex: Acid rain : $k = 7 \Rightarrow$ total $2^3 = 128$ possible nodel
complexity seerched bush model
 $M = 42$; $T = x_4$ (AL)
 $M = 53$; $(2) = 21$ x_1, x_3
 $(n_1 = 554$; $(3) = 35$ x_1, x_2, x_3, x_5, x_7
 $M = x_1$ (6) $x_1, x_2, x_3, x_3, x_5, x_7$
 $(n_1 = 28 (5) \cdot 1 - 24$ x_1, x_2, x_3, x_5, x_7

•

S

3) Now we need to choose between these kell moally found in 2). Which criterion should I use? (17]=p=kel

i)
$$R^2$$
 adj = $l - \frac{SSE}{N-l}$

ii) Mallows'
$$G = \frac{SSE}{\hat{\sigma}_{Form}^2} - n + 2|\pi|$$

 $VS SPSE = SSE + 2|\pi| \cdot \hat{\sigma}_{Form}^2$
iv) AIC = $n \cdot (m(\hat{\sigma}^2) + 2(|\pi|+1))$
 $\frac{SSE}{N-p_{H_{HT}}}$
iv) BIC = $n \cdot (m(\hat{\sigma}^2) + (m(n))(|\pi|+1))$

BIC gives non penally then AIC to large nodels.

4

Trensformation of response and predictors night improve the fit of the regression model.

The BoxCox trensform

$$g_{\lambda}(\lambda) = \begin{cases} y^{\lambda-1} & \lambda \neq 0 \\ \chi & (\Lambda(y) & \Lambda \neq 0 \end{cases}$$

Class of function

For Y= XB+E, En N(0,0° F) the best value of A on based on maximizing the likelihood proble $l(A) = -\frac{n}{2} \ln \left(\frac{SSEA}{n} \right) - (A-1) \sum_{i=1}^{n} \ln^{i} i$ SSEA (s the SSE when $g_A(Y)$ is the response R: boxcox(fit), see plot.

TMA4267 Linear statistical models

Part 2: Linear regression

February 20, 2017

Normal equations

$$\mathbf{Y} = \mathbf{X}\mathbf{\beta} + \mathbf{\epsilon}$$
 where $E(\mathbf{\epsilon}) = \mathbf{0}$ and $Cov(\mathbf{\epsilon}) = \sigma^2 \mathbf{I}$

Which of the following are *the normal equations*?

A
$$X\hat{\beta} = HY$$

B $\hat{\beta} = X(X^TX)^{-1}X^TY$
C $(X^TX)\hat{\beta} = X^TY$
D $(X^TX)Y = X^T\hat{\beta}$

The hat matrix

Design matrix **X** has *n* rows and *p* linearly independent columns. $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is called the hat-matrix.

Which of the following statements are NOT true?

A $H = H^{T} = H^{2}$ **B** rank(H) = p**C** HY = Y **D** H(I - H) = 0



$Y = X\beta + \varepsilon \text{ where } E(\varepsilon) = 0 \text{ and } Cov(\varepsilon) = \sigma^2 I$ $H = X(X^T X)^{-1} X^T$

An unbiased estimator for σ^2 is:

- **A** SSE/n **B** $\mathbf{Y}^T (\mathbf{I} \mathbf{H}) \mathbf{Y} / (n p)$
- **C** $(X^T X)^{-1} Y / (n p)$ **D** $(X^T X)^{-1} SSE / n$

Inference about β

$$\mathbf{Y} = \mathbf{X}\mathbf{\beta} + \mathbf{\varepsilon}$$
 where $\mathbf{\varepsilon} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$
and $\hat{\mathbf{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$.

What are the properties of $\hat{\beta}$?

- A Chi-squared distributed with n - pdegrees of freedom.
- C Multivariate normal with covariance matrix $(I H)\sigma^2$.
- B Chi-squared distributed with p degrees of freedom.
- **D** Multivariate normal with covariance matrix $(\mathbf{X}^T \mathbf{X})^{-1} \sigma^2$.

Happiness=money+sex+love+work

	Estimate	Std. Error	t value	Pr(> t)
money	0.009578	0.005213	1.837	0.0749
sex	-0.149008	0.418525	-0.356	0.7240
love	1.919279	0.295451	6.496	1.97e-07
work	0.476079	0.199389	2.388	0.0227

Which of the regression coefficient estimates has the largest estimated variance?

A moneyB sexC loveD work

Happiness=money+sex+love+work

The R^2 for the happiness-regression model is 71%. What does that mean?

- A The regression is significant for significance level 71%
- **B** The regression explains 71% of the variability in the data
- **c** The estimate for the variance σ^2 is 0.71
- D The covariates have a correlation of 0.71

Happiness

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.072081	0.852543	-0.085	0.9331
money	0.009578	0.005213	1.837	0.0749
sex	-0.149008	0.418525	-0.356	0.7240
love	1.919279	0.295451	6.496	1.97e-07
work	0.476079	0.199389	2.388	0.0227

For which β_j would we reject the null hypothesis $\beta_j = 0$ at significance level 1%?

A money B sex

C love D work

Best model



Which model does the BIC criterion report to be the best?

- A love+work B love
- **C** money+love+work **D** money+sex+love+work

What is this plot used for?



A Check residuals

C Assess linearity

- B Assess normality of residuals
- D Find transform of response
Correct?

Are you sure you want to read the correct answers? Maybe try first? The answers are explained on the next two slides.

Answers

- 1. C: The normal equation $(\mathbf{X}^{T}\mathbf{X})\hat{\boldsymbol{\beta}} = \mathbf{X}^{T}\mathbf{Y}$ is before you solve for $\hat{\boldsymbol{\beta}}$.
- 2. C: The hat matrix is symmetric and idempotent (so A is ok), and has rank p, but the reason for the name of the hat matrix is that is puts the hat on the Y so $HY = \hat{Y}$. We know that for symmetric projection matrices the two matrices H and (I H) are orthogonal so the product must be zero.

Answers

- 3. B: Since SSE has mean $(n p)\sigma^2$, then SSE/(n-p) must be an unbiased estimator for σ^2 . We know that (I - H) projects onto the space othogonal to the column space of the designmatrix, so that must have to do with SSE.
- 4. D: We know that linear combinations of multivariate normal random vectors are also multivariate normal (so the chisquare is not suitable). The residuals have (I H) as part of their covariance matrix, but $\hat{\beta}$ has not.

Answers

- 5. B: Sex has the largest estimated variance for regression estimate.
- 6. B: R^2 gives the percent of variability explained.
- 7. C: only love is significant on level 1%, since this is the only *p*-value below 0.01 (last column).
- 8. A: love+work has smallest BIC.
- 9. D: Box-Cox plot used to find transformation of response.